



Modulation of Brain Activity by Selective Attention to Audiovisual Dialogues

Alina Leminen^{1,2,3,4*}, Maxime Verwoert¹, Mona Moisala¹, Viljami Salmela¹, Patrik Wikman¹ and Kimmo Alho^{1,5*}

¹ Department of Psychology and Logopedics, Faculty of Medicine, University of Helsinki, Helsinki, Finland, ² Cognitive Science, Department of Digital Humanities, Helsinki Centre for Digital Humanities (Heldig), University of Helsinki, Helsinki, Finland, ³ Cognitive Brain Research Unit, Department of Psychology and Logopedics, Faculty of Medicine, University of Helsinki, Helsinki, Finland, ⁴ Center for Cognition and Decision Making, Institute of Cognitive Neuroscience, National Research University – Higher School of Economics, Moscow, Russia, ⁵ Advanced Magnetic Imaging Centre, Aalto Neuroimaging, Aalto University, Espoo, Finland

OPEN ACCESS

Edited by:

Claude Alain,
Rotman Research Institute (RRI),
Canada

Reviewed by:

Michael S. Beauchamp,
Baylor College of Medicine,
United States
Yang Zhang,
University of Minnesota Twin Cities,
United States

*Correspondence:

Alina Leminen
alina.leminen@helsinki.fi
Kimmo Alho
kimmo.alho@helsinki.fi

Specialty section:

This article was submitted to
Auditory Cognitive Neuroscience,
a section of the journal
Frontiers in Neuroscience

Received: 24 September 2019

Accepted: 09 April 2020

Published: 12 May 2020

Citation:

Leminen A, Verwoert M,
Moisala M, Salmela V, Wikman P and
Alho K (2020) Modulation of Brain
Activity by Selective Attention
to Audiovisual Dialogues.
Front. Neurosci. 14:436.
doi: 10.3389/fnins.2020.00436

In real-life noisy situations, we can selectively attend to conversations in the presence of irrelevant voices, but neurocognitive mechanisms in such natural listening situations remain largely unexplored. Previous research has shown distributed activity in the mid superior temporal gyrus (STG) and sulcus (STS) while listening to speech and human voices, in the posterior STS and fusiform gyrus when combining auditory, visual and linguistic information, as well as in left-hemisphere temporal and frontal cortical areas during comprehension. In the present functional magnetic resonance imaging (fMRI) study, we investigated how selective attention modulates neural responses to naturalistic audiovisual dialogues. Our healthy adult participants ($N = 15$) selectively attended to video-taped dialogues between a man and woman in the presence of irrelevant continuous speech in the background. We modulated the auditory quality of dialogues with noise vocoding and their visual quality by masking speech-related facial movements. Both increased auditory quality and increased visual quality were associated with bilateral activity enhancements in the STG/STS. In addition, decreased audiovisual stimulus quality elicited enhanced fronto-parietal activity, presumably reflecting increased attentional demands. Finally, attention to the dialogues, in relation to a control task where a fixation cross was attended and the dialogue ignored, yielded enhanced activity in the left planum polare, angular gyrus, the right temporal pole, as well as in the orbitofrontal/ventromedial prefrontal cortex and posterior cingulate gyrus. Our findings suggest that naturalistic conversations effectively engage participants and reveal brain networks related to social perception in addition to speech and semantic processing networks.

Keywords: selective attention, noise vocoding, audiovisual integration, social perception, speech, visual speech, fMRI

INTRODUCTION

In everyday life, we are often faced with multiple speaker situations, for instance, when dining in a crowded restaurant or talking to a friend while hearing a radio in the background. Such situations require segregation of speech streams originating from different sources and selection of one of

the streams for further processing. The neural mechanisms through which this type of attentional selection is achieved are not yet fully understood (e.g., Rimmele et al., 2015).

A meta-analysis (Alho et al., 2014) of functional magnetic resonance imaging (fMRI) studies on stimulus-dependent sound processing and attention-related modulations in the auditory cortex showed that speech and voice processing activate overlapping areas in the mid superior temporal gyrus and sulcus bilaterally (STG and STS, respectively). Furthermore, selective attention to continuous speech appeared to modulate activity predominantly in the same areas (Alho et al., 2014). Importantly, selectively attending to a particular speaker in a multi-talker situation results in the STG activity that represents the spectral and temporal features of attended speech, as if participants were listening only to that speech stream (Mesgarani and Chang, 2012). In other words, the human auditory system restores the representation of an attended speaker while suppressing irrelevant or competing speech.

In addition to STG/STS, selective attention to non-speech sounds engages prefrontal and parietal cortical areas (Tzourio et al., 1997; Alho et al., 1999; Zatorre et al., 1999; Degerman et al., 2006), which has been associated with top-down control needed to select attended sounds and reject irrelevant sounds. Selective attention to continuous speech, however, does not appear to markedly engage prefrontal and superior parietal areas (Alho et al., 2003, 2006; Scott et al., 2004). This is most probably because selective listening to speech is a highly automatized process, less dependent on fronto-parietal attentional control (Alho et al., 2006; see also Mesgarani and Chang, 2012). Such automaticity might be due to listeners' lifelong experience in listening to speech. However, initial orienting of attention to one of three concurrent speech streams has yielded enhanced activation in the fronto-parietal network, hence, purportedly engaging an attentional top-down control mechanism (Alho et al., 2015; Hill and Miller, 2010).

Natural situations with multiple speakers might not only be complicated by the demand to listen selectively to one speech stream while ignoring competing speech, but also by degraded quality of the attended speech (e.g., when talking in a noisy café on the phone with a poor signal). Studies addressing the comprehension of degraded (e.g., noise-vocoded) speech involving only one speech stream have reported increased activity in the posterior parietal cortex (Obleser et al., 2007) and frontal operculum (Davis and Johnsrude, 2003), as compared to more intelligible speech. Listening to degraded, yet intelligible and highly predictable speech, in turn, elicits activity in the dorsolateral prefrontal cortex, posterior cingulate cortex, and angular gyrus (e.g., Obleser et al., 2007). Moreover, the amount of spectral detail in speech signal was found to correlate with STS and left inferior frontal gyrus (IFG) activity, regardless of semantic predictability (Obleser et al., 2007). McGettigan et al. (2012) observed increasing activity along the length of left dorsolateral temporal cortex, in the right dorsolateral prefrontal cortex and bilateral IFG, but decreasing activation in the middle cingulate, middle frontal, inferior occipital, and parietal cortices associated with increasing auditory quality. Listening to degraded

speech has also activated the left IFG, attributed to higher-order linguistic comprehension (Davis and Johnsrude, 2003) and the dorsal fronto-parietal network, related to top-down control of attention (Obleser et al., 2007). Overall, increased speech intelligibility enhances activity in the STS (Scott et al., 2000; Obleser et al., 2007; McGettigan et al., 2012), STG (Davis and Johnsrude, 2003), middle temporal gyrus (MTG; Davis and Johnsrude, 2003), and left IFG (Davis and Johnsrude, 2003; Obleser et al., 2007; McGettigan et al., 2012). Increased activity in these areas may be related to enhanced speech comprehension due to increasing availability of linguistic information.

The studies described above, however, used only single-speaker paradigms. Evans et al. (2016) examined how different masking sounds are processed in the human brain. They used a selective attention paradigm with two speech streams, namely, a masked stream and a target stream. The target speech was always clear, whilst the masked speech was either clear, spectrally rotated or noise-modulated. Increased intelligibility of the masked speech activated the left posterior STG/STS, however, less extensively than a clear single speech alone. This was taken to suggest that syntactic and other higher order properties of masking speech are not actively processed and the masker sounds may be actively suppressed already at early processing stages (see also Mesgarani and Chang, 2012). In contrast, the masked speech yielded increased activation in the frontal (bilateral middle frontal gyrus, left superior orbital gyrus, right IFG), parietal (left inferior and superior parietal lobule) and middle/anterior cingulate cortices, as well as in the frontal operculum and insula. These activations were suggested to reflect increased attentional and control processes. The results corroborate those from earlier positron emission tomography (PET) studies (e.g., Scott et al., 2004) on selective attention to a target speaker in the presence of another speaker (speech-in-speech) or noise (speech-in-noise). More specifically, Scott et al. (2004) found more activity in the bilateral STG for speech-in-speech than speech-in-noise, whereas speech-in-noise elicited more activity in the left prefrontal and right parietal cortex than speech-in-speech. Scott and colleagues suggested that these additional areas might be engaged to facilitate speech comprehension or that they are related to top-down attentional control. Correspondingly, Wild et al. (2012) reported activations in frontal areas (including the left IFG) that were only present when the participants selectively attended to the target speech among non-speech distractors. In contrast to studies reporting increased left IFG activations to increased intelligibility of degraded speech (Davis and Johnsrude, 2003; Obleser et al., 2007; McGettigan et al., 2012), Wild et al. (2012) found greater activity in the left IFG for degraded than for clear target speech. By contrast, STS activity was increased with decreasing speech intelligibility, regardless of attention. Increased activity for attended degraded speech was proposed to reflect "the improvement in intelligibility afforded by explicit, effortful processing, or by additional cognitive processes (such as perceptual learning) that are engaged under directed attention" (Wild et al., 2012, p. 14019). The authors further suggested that top-down influences on early auditory processing might facilitate speech comprehension in difficult listening situations.

The majority of fMRI studies on selective attention to speech have used only auditory speech stimuli (e.g., Alho et al., 2003, 2006; Wild et al., 2012; Evans et al., 2016; Puschmann et al., 2017). However, natural conversations often include also visual speech information. Integrating a voice with mouth movements (i.e., visual speech) facilitates speech understanding in relation to mere listening (Sumbly and Pollack, 1954). In accordance, fMRI studies on listening to speech have shown that the presence of visual speech enhances activity in the auditory cortex and higher order speech-processing areas (e.g., Bishop and Miller, 2009; McGettigan et al., 2012). A related magnetoencephalography (MEG) study showed that the presence of visual speech enhances auditory-cortex activity that follows the temporal amplitude envelope of attended speech (Zion Golumbic et al., 2013; for similar electroencephalography (EEG) evidence, see O'Sullivan et al., 2015). Facilitation of speech comprehension by visual speech holds especially true for noisy situations (e.g., Sumbly and Pollack, 1954) and degraded quality of attended speech (e.g., McGettigan et al., 2012; Zion Golumbic et al., 2013). Some fMRI studies have suggested maximal facilitation of speech comprehension by visual speech at intermediate signal-to-noise ratios of auditory information (Ross et al., 2007; McGettigan et al., 2012).

Degraded speech increases demands for fronto-parietal top-down control (Davis and Johnsrude, 2003; Evans et al., 2016), whereas adding visual speech appears to facilitate selective attention (Sumbly and Pollack, 1954; Zion Golumbic et al., 2013). However, it is still unknown whether fronto-parietal areas are activated during selective attention to visually degraded speech. Moreover, an earlier study that employed a factorial design with different levels of auditory and visual clarity in sentences (McGettigan et al., 2012) did not include an unmodulated (clear) visual and auditory condition. Hence, to our knowledge, brain responses to continuous naturalistic dialogues with varying audio-visual speech quality have not been systematically examined before.

In the current study, we collected whole-head fMRI data in order to identify brain regions critical for selective attention to natural audiovisual speech. More specifically, we examined attention-related modulations in the auditory cortex and associated fronto-parietal activity during selective attention to audiovisual dialogues. In addition, we assessed an interplay between auditory and visual quality manipulations. We also included clear auditory and visual stimulus conditions to investigate brain areas activated during selective attention to naturalistic dialogues in the presence of irrelevant clear speech in the background. Our experimental setup might be regarded as mimicking watching a talk show on a TV while a radio program is playing on the background. Comparing brain activity during attention to the dialogues with activity during control conditions, where the dialogues are ignored and fixation cross is to be attended, allowed us to determine attention-related top-down effects and distinguish them from stimulus-dependent bottom-up effects (Alho et al., 2014).

We predicted that both increased speech intelligibility and increased amount of visual speech information in the attended speech would be associated with stronger stimulus-dependent

activity in the STG/STS as well as subsequent activity in brain areas involved in linguistic processing. Moreover, we hypothesized that degrading auditory or visual quality of attended speech might be related to increased fronto-parietal activity due to enhanced attentional demands. Finally, we were interested to see whether attention to audiovisual speech and the quality of this speech would have interactions in some brain areas involved in auditory, visual or linguistic processing, or in the control of attention.

METHODS

Participants

Fifteen healthy right-handed adult volunteers (5 males, age range 20–38 years, mean 25.3 years) participated in the present study. All participants were native Finnish speakers with normal hearing, normal or corrected-to-normal vision, and no history of psychiatric or neurological illnesses. Handedness was verified by the Edinburgh Handedness Inventory (Oldfield, 1971). An informed written consent was obtained from each participant before the experiment. The experimental protocol was approved by the Ethics Review Board in the Humanities and Social and Behavioral Sciences, University of Helsinki.

Stimuli

Stimulus Preparation

The stimuli consisted of 36 video clips showing scripted spoken dialogues (see **Table 1** for an example of a dialogue). The topics of dialogues were neutral, such as weather, vacation, and study plans. The syntactic structure of dialogues was matched as closely as possible. An independent native Finnish speaker subsequently verified the neutrality of dialogues as well as their meaningfulness and grammaticality. Each dialogue always consisted of seven lines spoken alternately by two actors, and each line contained 9–13 words.

The stimulus recordings took place in a soundproof studio. The video clips were recorded with a wide angle (23.5 mm G lens) HXR-NX70E digital video camera (SONY Corporation, Tokyo, Japan). Two external microphones were attached to the camera in order to record the left and right audio channels separately (48 kHz sampling frequency, 16-bit quantization).

The actors were two native Finnish speakers (a male and female university student recruited for the recording purposes). They were unaware of the experimental setup and were compensated for their work. The actors memorized the dialogues beforehand but uttered their lines with a natural pace. An external prompter (programmed with Matlab version R2016, Mathworks Inc., Natick, MA, United States) was used to remind each actor to hold a pause before uttering the next line. The pause duration information was used in the subsequent fMRI data processing. The mean duration of dialogues was 60 s (range 55–65 s) with mean line duration of 5.4 s and inter-line pause duration of 3.4 s. Half of the dialogues started with the female speaker and the other half with the male speaker. The speakers sat next to one another with their faces slightly tilted toward each other, making the visual speech setting as natural as possible while maintaining

visual speech information visible to a viewer. The video data were then edited with Corel VideoStudio Pro X 8 software (Corel Corporation, Ottawa, ON, Canada) and, finally, with Matlab, see below. The video clips were cut into separate dialogues with 720 ms (18 frames) before the first and after the last spoken words. Thereafter, the videos were split into separate video and audio channels for subsequent editing with Adobe Audition CS6 (Adobe Systems Inc., San Jose, CA, United States) software. The audio channels were then converted to mono, cleaned from all non-voice background sounds, low-pass filtered at 5000 Hz, and scaled to have the same peak sound energy in all dialogues.

In addition to the natural speech, the audio data were noise-vocoded (Shannon et al., 1995; Davis and Johnsrude, 2003) using Praat software (version 6.0.27; Boersma and Weenink, 2001). The audio files were divided into 2 and 4 logarithmically spaced frequency bands between 300 and 5000 Hz (2 band cut-off points: 300, 1385, 5000 and 4 band cut-off points 300, 684, 1385, 2665, 5000). The filter bandwidths were set to represent equal distances along the basilar membrane (according to the Greenwood (1990) equation relating filter position to best frequency). The amplitude envelope from each frequency band was extracted using the standard Praat algorithm. The extracted envelope was then applied to band-pass filtered noise in the same frequency bands. Then, the resulting bands of modulated noise were recombined to produce the distorted speech. Noise vocoded speech sounds like a harsh robotic whisper (Davis and Johnsrude, 2003). Finally, the unchanged F0 (frequencies 0–300 Hz) was added to the noise-vocoded speech in order to maintain the speakers' gender identity clearly perceivable and their voices distinguishable from the irrelevant voice speaking in the background (see below). The speech was perceived to be hardly intelligible with 3 frequency bands (i.e., 2 noise-vocoded bands and the intact F0 band) and quite intelligible with 5 frequency bands (i.e., 4 noise-vocoded bands and the intact F0 band). These two frequency-band manipulations for noise-vocoding were assumed to be optimal

for our study on the basis of a behavioral pilot experiment. In this pilot experiment, 5 listeners (not included in the actual fMRI experiment) rated the intelligibility of seven dialogues noise-vocoded across a wide range of frequency bands (2, 4, 6, 8, 10, 12 or 16) with a non-vocoded F0 band. The participants listened to the dialogues one line at a time and provided a typed report on what they could hear. On average, for 2 and 4 noise-vocoded bands, 26.2% ($SD = 18.6\%$) and 76.4% ($SD = 10.3\%$) of the lines were perceived correctly.

In addition to manipulating auditory information, we parametrically varied the amount of visual speech seen by the participants. This was done by adding different amounts of dynamic white noise onto the region in the videos showing the speakers' faces. Noise was added with Matlab R2016 using built-in functions and custom-made scripts with the following procedure. First, we constructed Gaussian masks for both faces (faces localized from the first frame of the video with Matlab's `vision.CascadeObjectDetector`). Then we generated two samples (one for each face) of white noise (using Matlab's `randn` function), and multiplied the noises with the facemasks in order to add noise smoothly only onto the faces. The same sample of noise was added to R, G, and B channels. This was repeated for every frame, and thus the noise was dynamic and it changed in every frame. To get different levels of visual quality, the amount of added noise was scaled so that the root-mean-contrast of the low-quality videos were 20 and 15% lower than the contrast of the highest quality video. Five experienced viewers (not included in the actual fMRI experiment) confirmed that adding the noise reduced the visual quality so that the mouth movements and facial features were only poorly visible at highest noise level (Figure 1).

In the final step of stimulus preparation procedure, we recombined the "poor," "medium," and "good" auditory quality sound files (with 2 noise-vocoded bands and an intact F0 band, 4 noise-vocoded bands and an intact F0 band, and clear intact speech, respectively) with the "poor," "medium," and "good" visual quality video files (more masked poorly perceivable visual speech, less masked quite perceivable visual speech, and unmasked clear visual speech, respectively) video files using a custom-made Matlab script. The resulting videos were then compressed using VirtualDub software¹. Example stimuli of the three experimental conditions (good auditory quality and good visual quality; medium auditory quality and medium visual quality; poor auditory quality and poor visual quality) can be found online². Written informed consent was obtained from the actors to publish identifiable image information.

Taken together, each dialogue had 3 visual and 3 auditory quality variants, which resulted in altogether 9 experimental conditions, one for each quality combination (e.g., poor visual and good auditory quality) with three dialogues in each. All combinations were presented to the participants but each participant saw a different variant of each dialogue.

Furthermore, to increase the attentional load, we added continuous background speech as an auditory distractor. For this purpose, we chose a cultural history audio book (the Finnish

TABLE 1 | Example of one natural speech dialogue by two actors (A and B) used in the experiment.

| Dialogue lines | Approximate english translation |
|---|--|
| A: Pitäisi kohta käydä kaupassa hakemassa välipalaa. Ostanko sinullekin jotain? | A: I should go soon to the store and get something to eat. Should I get something for you as well? |
| B: Ei kiitos tarvitse, minä pakkasin leipää ja jogurttia tänä aamuna lounaaksi. | B: No thanks, I packed bread and yogurt with me for lunch today. |
| A: Hyvä on, mutta haluatko tulla mukaan seuraksi kauppaan kuitenkin? | A: Okay, but would you still like to come along with me to the store? |
| B: Mielelläni, voisin katsoa, jos löytäisin sieltä jotain syötävää myöhemmälle. | B: With pleasure, I could see if I would find something to eat later. |
| A: Haluatko tulla kanssani puistoon syömään, kun olemme tulleet kaupasta? | A: Would you like to come with me to the park to eat after visiting the store? |
| B: Ulkona on aika kylmä tänään. Mentäisiinkö mieluummin jonnekin sisälle? | B: It's quite cold outside today. Should we rather go somewhere inside? |
| A: Totta, voisimme siinä tapauksessa syödä täällä yliopiston kahvihuoneessa. | A: True, in that case we could eat here at the university in the coffee room. |

¹<http://www.virtualdub.org>

²<https://osf.io/h9er7/>

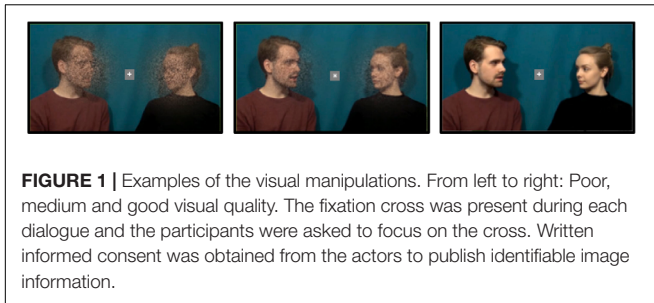


FIGURE 1 | Examples of the visual manipulations. From left to right: Poor, medium and good visual quality. The fixation cross was present during each dialogue and the participants were asked to focus on the cross. Written informed consent was obtained from the actors to publish identifiable image information.

translation of *The Autumn of the Middle Ages* by Johan Huizinga), which is freely distributed online by the Finnish Broadcasting Company (Yleisradio, YLE)³. The book was read by a female professional Finnish-native actor. In order for the F0 in this auditory distractor to be perceived approximately equidistant from the F0s of our female (200 Hz) and male (122 Hz) actors, we manipulated the F0 of the reader's voice by using square root of the mean of the female and male voices in the recorded video clips. After some further manipulations based on the estimation of three experienced listeners, the resulting F0 was 156 Hz. The F0 manipulation was performed in Audacity software⁴. The background speech was otherwise presented in its natural form and low-pass filtered at 5000 Hz to match the audio used in the experimental conditions. The audiobook was always presented as clear (i.e., non-vocoded) speech in the background. In addition, loudness differences between attended and unattended speech were kept minimal, as verified by three experienced listeners.

Procedure

Stimulus presentation was controlled through a script written in Presentation 20.0 software (Neurobehavioral Systems Inc., Berkeley, CA, United States). The video clips were projected onto a mirror mounted on the head coil and presented in the middle of the screen. All auditory stimuli were presented binaurally through insert earphones (Sensimetrics model S14; Sensimetrics, Malden, MA, United States). The experiment consisted of 3 functional runs with all 9 experimental conditions (Auditory Quality either poor, medium or good and Visual Quality either poor, medium or good) presented in each run along with 2 visual control conditions. The order of conditions was also randomized; however, the visual control conditions were always presented at the 6th and 7th place within a run. There was a small break of 40 s between these two dialogues. During the rest period, the participants were asked to focus on the fixation cross. Within all three functional runs, the order of the conditions was randomized for each participant. The competing audio distractor (audiobook) was presented 500–2000 ms before video onset and stopped at the offset of the video. The differing durations of dialogues were compensated for by inserting periods with a fixation cross between the instruction and the onset of the dialogue, keeping the overall trial durations constant. The intensity of the sounds was individually set to a loud, but pleasant

³<https://areena.yle.fi/1-3529001>

⁴<https://sourceforge.net/projects/audacity/>

level, and was approximately 80 dB SPL as measured from the tip of the earphones.

Attention-to-Speech Conditions

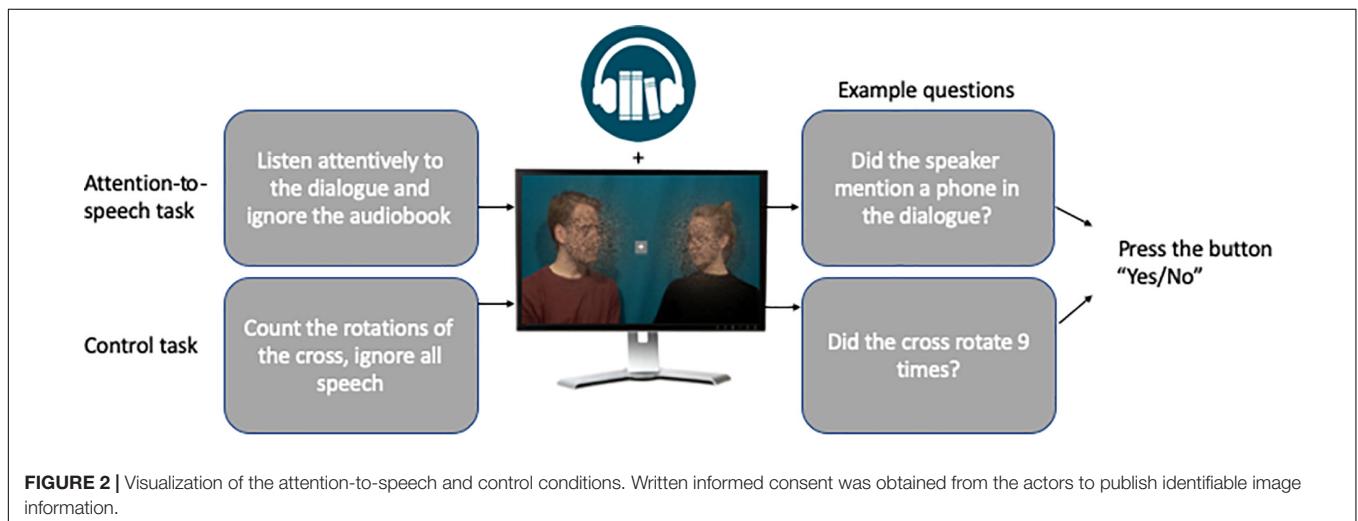
In the attention-to-speech conditions, the participants were asked to attentively watch the videos, ignore the background speech, and after each 7-line dialogue answer to seven questions, one question related to each line of the dialogue. More specifically, they were instructed to answer whether a certain topic was discussed in a particular line (see **Table 2**) by pressing the “Yes” or “No” button on a response pad (LUMItouch, Photon control Inc., Richmond, BC, Canada) with their right index or middle finger, respectively. Each written question was presented on the screen for 2 s, during which the participant gave his/her answer. Regardless of the duration of the participant's answer, the next question always started 2 s after the previous one. After the 7 questions, the participants were provided with immediate feedback on their performance (number of the correct answers), and the fixed duration of feedback was 2 s. The next dialogue was presented after a short (2-s) written instruction, telling whether the task was an attention-to-speech or a control (see section Attention-to-the-Fixation-Cross Condition) dialogue, followed by the fixation cross period, presented for 3–13 s to make all the trials equally long. All video clips had a rotating white fixation cross (inside a light gray box), placed in the middle of the screen, with a minimum of 9 and maximum of 15 rotations at a random interval but with minimum of 3 s between rotations. In the attention-to-speech condition, the participants' task was to ignore the cross and concentrate on viewing the people speaking.

Attention-to-the-Fixation-Cross Condition

In addition to the attention-to-speech conditions, we included two control conditions. These consisted of videos with a combination of good auditory and good visual quality and a combination of poor auditory and poor visual quality. Note that only these auditory and visual quality combinations were included into control conditions, because adding the other seven combinations to all three runs would have prolonged the total duration of the experiment by about 30 min and made the experiment too long to be conducted in a single experimental session. The dialogues in these control conditions were the ones not used in the attention-to-speech conditions. Identically to the attention-to-speech conditions, all video clips had a rotating white fixation cross (inside a light gray box), placed in the middle of the screen, with a minimum of 9 and maximum of 15 rotations from “×” to “+”, or vice versa, at random intervals with a minimum of 3 s between the rotations. Thus, the attention-to-the-fixation-cross conditions used the same setup as the attention-to-speech conditions, but the task of the participants was to concentrate on counting the number of times the fixation cross rotated and ignore the dialogue and the background voice. After each control block, the participants were presented with seven questions (“Did the cross turn X times?”; the X being 9, 10, 11, 12, 13, 14, or 15 in an ascending order), and they were asked to answer each question with “Yes” or “No” by pressing the corresponding button on the response pad with their right index or middle

TABLE 2 | Example of quiz questions of a practice dialogue.

| Dialogue | English translation | A related quiz question "Did the speakers discuss this topic?" | Correct answer |
|---|---|---|----------------|
| A: Ostin uuden puhelimen ja siinä on niin paljon toimintoja, että olen sen kanssa ihan hukassa. | A: I bought a new phone and it has so many features that I am completely lost. | Puhuja hukkasi puhelimensa./The speaker lost his/her phone. | No |
| B: Ai niin joo, sinulla oli ennen sellainen ikivanha kännykkä, joka ei ollut edes älypuhelin. | B: Oh yes, you used to have that ancient phone, which was not even a smart phone. | Puhujan kännykkä oli vanha./The speaker's phone was old. | Yes |
| A: Joo, ja se oli aivan hyvä puhelin, siihen asti kunnes kissa pudotti sen pöydältä lattialle, se oli sitten siinä. | A: Yes, and it was a perfectly good phone until my cat dropped it on the floor from a table, and that was it. | Koira rikkoi puhelimen./The dog broke the phone. | No |
| B: Minä kun luulin, että vanhat kännykät kestävät kaiken eivätkä menisi mistään rikki. | B: I thought that old phones take all hits and wouldn't break at all. | Puhuja ihaili uutta puhelinta./The speaker was admiring the new phone. | No |
| A: No se on kyllä pudonnut monta kertaa, mutta kestänyt kaikki iskut mutta tämä taisi olla sille liikaa. | A: It has indeed fallen many times and always stayed intact but now this was too much for it. | Vanha kännykkä kesti iskut./The old mobile endured all hits. | Yes |
| B: No, mutta toivotaan, että tässä uudessa kännykässäsi kestää akku hyvin ja olet siihen muutenkin tyytyväinen. | B: Well, let's hope that your new phone has a long-lasting battery and that you are satisfied with it in all aspects. | Puhujalla ei ollut laturia mukanaan./The speaker did not have a charger with him/her. | No |
| A: Nyt on vielä vähän hankalaa, enkä osaa sitä oikein käyttää mutta kyllä se varmaan tästä! | A: It is still a bit difficult and I really don't know how to use it but I think it will be fine! | Puhuja ei osannut käyttää älypuhelin./The speaker did not know how to use the smartphone. | Yes |



finger, respectively. Thereafter, the participants were provided with immediate feedback (e.g., “6/7 correct”). For a schematic presentation of one trial, see **Figure 2**.

Practice Trial

Before the actual fMRI scanning, the participants were familiarized with the task outside the scanner by viewing one practice dialogue with all conditions and answering questions related to its content.

Data Acquisition

Functional brain imaging was carried out with 3T MAGNETOM Skyra whole-body scanner (Siemens Healthcare, Erlangen, Germany) using a 30-channel head coil. The functional echo planar (EPI) images were acquired with an imaging area consisting of 43 contiguous oblique axial slices (TR 2530 ms, TE 32 ms, flip angle 75°, voxel matrix 64 × 64, field of view 20 cm, slice thickness 3.0 mm, in-plane resolution

3.1 mm × 3.1 mm × 3.0 mm). Three functional runs of 368 volumes were measured for each participant. A total of 1158 functional volumes were obtained in one session (session duration approximately 50 min). High-resolution anatomical images (voxel matrix 256 × 256, in-plane resolution 1 mm × 1 mm × 1 mm) were acquired from each participant prior to the functional runs.

Data Analysis

The fMRI data were pre-processed and analyzed in Statistical Parametric Mapping (SPM12; Wellcome Trust Centre for Neuroimaging, London, United Kingdom). The first 4 volumes in each run were dummies and were discarded in further analysis of the data, leaving 382 total volumes per run to be analyzed. The data were slice-time corrected, motion corrected, realigned to the middle image from each run, high-pass filtered (cutoff 1/260 Hz) and spatially smoothed with a Gaussian kernel of 6 mm. The images were normalized to MNI space using a standard

pre-processing function in Conn software (Whitfield-Gabrieli and Nieto-Castanon, 2012). For the first-level statistical analysis, the general linear model was created including a regressor for each condition. Separate regressors were also included for (1) the instructions and the responses from the participant and (2) the quiz. This resulted in 13 regressors in total. Additionally, six movement parameters (3 translations, 3 rotations) were included as nuisance regressors. The conditions were modeled using a standard boxcar model. For the second-level analysis, we used the Multivariate and Repeated Measures (MRM) toolbox (McFarquhar et al., 2016). The contrast images of the nine experimental conditions compared to rest from each participant were entered into a 3×3 full factorial repeated-measures analysis of variance (ANOVA) with factors Visual Quality (3 levels: poor, medium, good) and Auditory Quality (3 levels: poor, medium, good). Within this model, F-contrasts were computed for the main effects and the interaction effect. A separate 2×2 repeated-measures ANOVA was conducted to account for stimulus quality and attentional effects. This additional ANOVA included factors Audiovisual Quality (2 levels: poor auditory and poor visual quality vs. good auditory and good visual quality) and Attention (2 levels: attention to speech vs. attention to the fixation cross). All reported contrasts were thresholded voxel-wise at $p < 0.001$ with a cluster extend threshold of 100 voxels, resulting activity maps that were family-wise error (FWE) corrected at the cluster level, $p(\text{FWE}) < 0.05$.

Statistical analyses of the performance data, that is, responses to the quiz questions during the attention-to-speech condition were submitted to the repeated-measures ANOVA with factors Visual Quality (poor, medium, good) and Auditory Quality (poor, medium, good). Responses to the quiz questions during the attention-to-the-fixation cross condition were submitted to the two-tailed pairwise t -test. For all analyses, the Greenhouse-Geisser correction was applied when the assumption of sphericity was violated. IBM SPSS Statistics 24 (IBM SPSS, Armonk, NY, United States) was used for conducting these analyses.

In addition to brain areas showing significant effects of factors included in ANOVAs, we studied brain activity in regions of interest (ROIs) known to be involved in low-level auditory processing and speech processing. These ROIs were located bilaterally in Heschl's gyrus (HG), the anterior, mid and posterior STG, as well as in Broca's area in the left hemisphere and its right hemisphere analogue (Liakakis et al., 2011; Alho et al., 2014; Liebenthal et al., 2014). In addition, due to our focus on visual speech processing, there were additional ROIs in the left and right fusiform face area (FFA; Grill-Spector and Weiner, 2014). These ROIs were based on the Harvard and Oxford cortical structural atlas⁵.

RESULTS

Behavioral Results

The mean performance scores for the attention-to-speech conditions are shown in **Figure 3**. Behavioral results demonstrated a significant main effect of Auditory Quality

$[F(2, 28) = 57.57, p = 0.001, \eta_p^2 = 0.80]$ and a significant main effect of Visual Quality $[F(2, 28) = 8.2, p = 0.002, \eta_p^2 = 0.37]$. Although visual quality appeared to have a slightly stronger effect on performance when the auditory quality was poor than when it was medium or good (see **Figure 5**), the interaction between the two factors did not reach significance $[F(4, 56) = 1.64, p = 0.176]$. Bonferroni-corrected *post hoc* tests for Auditory Quality revealed significant differences between all Auditory Quality conditions (for all comparisons, $p < 0.001$). *Post hoc* tests for Visual Quality revealed significant differences between the poor and good quality conditions ($p < 0.001$) and between medium and good quality conditions ($p = 0.029$).

The mean performance scores for the attention-to-the-fixation-cross conditions were 6.7/7 (SEM 0.16/7) and 6.5/7 (SEM 0.36/7) for the poor audiovisual quality condition and good audiovisual quality condition, respectively ($p = 0.078$).

fMRI Results

Auditory and Visual Quality

Figure 4 depicts the brain areas where the 3×3 ANOVA showed significant main effects of Auditory Quality (3 levels) and **Figure 5** significant main effects of Visual Quality (3 levels) on brain activity measured during attention to speech. **Figures 4** and **5** also depict mean parameter estimates of significant clusters, displaying the direction of the observed cluster effect. No significant interactions between these factors were found with the applied significance threshold.

As seen in **Figure 4A**, Auditory Quality showed a significant effect on brain activity in the STG/STS bilaterally, these effects extending from mid-STG/STS areas to the temporal poles, the left angular gyrus, and the left superior frontal gyrus. **Figure 4B** demonstrates that in all these areas activity was enhanced with increasing auditory quality. **Figures 4C,D** also depict activity in additional ROIs bilaterally: HG, the anterior, mid and posterior STG, and Broca's area in the left hemisphere and its right-hemisphere analogue.

Visual Quality, in turn, had a significant effect on brain activity in the temporal and occipital cortices (**Figure 5A**). As seen in **Figure 5B**, increasing visual quality was associated with enhanced activity in the STS bilaterally, this activity extended in the right hemisphere even to the temporal pole, and in the left hemisphere to the inferior frontal gyrus. However, **Figure 5C** shows that activations increased with decreasing visual quality in the right fusiform gyrus and in the bilateral middle occipital gyrus predominantly in the left hemisphere. **Figure 5D** depicts activity in the bilateral FFA ROIs.

Additional ANOVAs conducted separately for activity in the HG, Broca's area and FFA ROIs included the factors Auditory Quality, Visual Quality, and Hemisphere (to avoid double-dipping, no such ANOVAs were performed for the STG ROIs covered already by the aforementioned 3×3 ANOVA). The results indicated a significant main effect of Hemisphere for all ROIs (HG $[F(1, 14) = 19.09, p = 0.001, \eta_p^2 = 0.58]$, Broca's area $[F(1, 14) = 9.53, p = 0.008, \eta_p^2 = 0.41]$ and FFA $[F(1, 14) = 11.81, p = 0.004, \eta_p^2 = 0.46]$). No other main or interaction effects were found for the HG and FFA ROIs, however, for the Broca's area ROI there was a significant interaction effect between

⁵<https://identifiers.org/neurovault.collection:262>

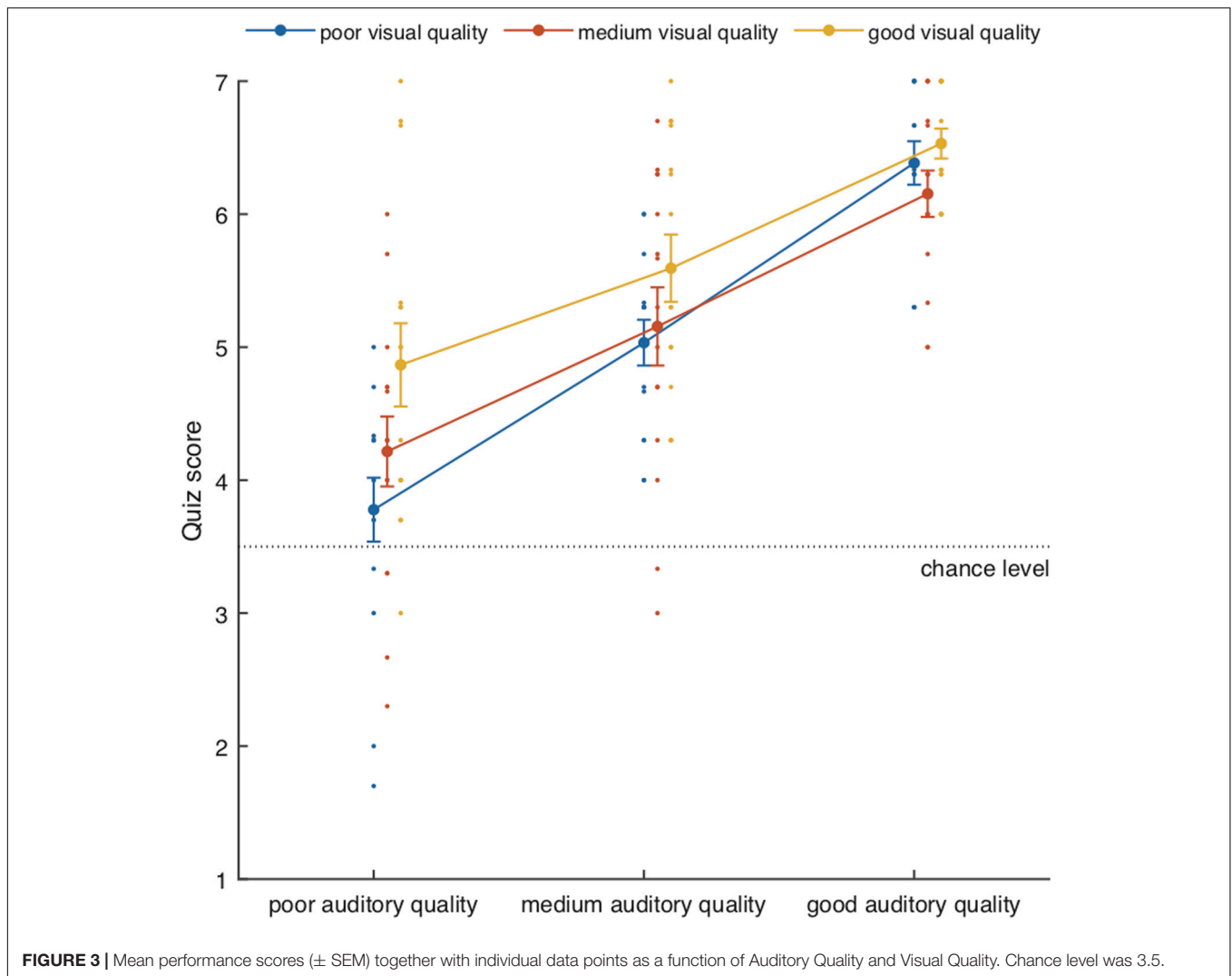


FIGURE 3 | Mean performance scores (\pm SEM) together with individual data points as a function of Auditory Quality and Visual Quality. Chance level was 3.5.

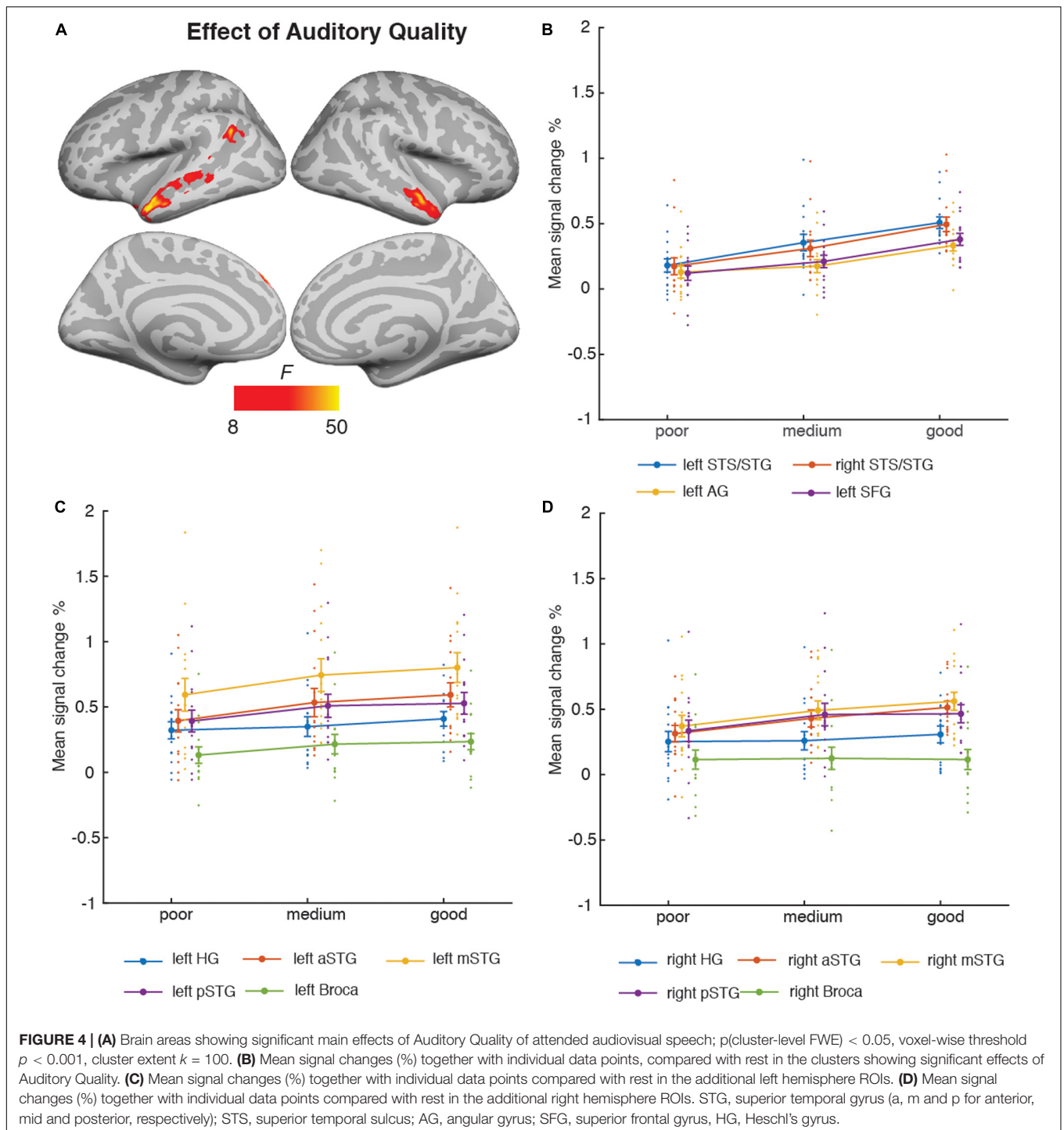
Hemisphere and Auditory Quality [$F(2, 28) = 7.4, p = 0.007, \eta_p^2 = 0.35$].

Attention-to-Speech vs. Attention-to-the-Fixation Cross

Figure 6 depicts the brain areas where the 2×2 ANOVA showed significant main effects of Audiovisual Quality (2 levels: poor auditory and poor visual quality vs. good auditory and good visual quality). **Figure 7** shows the brain areas where the 2×2 ANOVA showed significant main effects of Attention (2 levels: attention-to-speech vs. attention-to-the-fixation cross) on brain activity. **Figures 6** and **7** also depict mean parameter estimates of significant clusters, displaying the direction of the observed cluster effect. No significant interactions between the factors Attention and Audiovisual Quality were observed with the applied significance threshold. As seen in **Figure 6A**, in the 2×2 ANOVA, there was a significant effect of Audiovisual Quality bilaterally in the STG/STS, these activations extended to the temporal poles, as well as to the left superior parietal lobule, left precuneus, the dorsal part of the right inferior

parietal lobule, and in the middle occipital gyrus bilaterally. As seen in **Figure 6B**, in the left and right superior temporal gyri, activity was enhanced with increasing audiovisual quality both during attention-to-speech and attention-to-the-fixation cross. In contrast, in both attention conditions, activity was higher in the left superior parietal lobule, the right inferior parietal, the left precuneus, and bilateral middle occipital gyrus for poorer audiovisual quality (**Figure 6C**). **Figure 6C** shows activity in additional ROIs: the left HG and anterior, mid and, posterior STG, Broca's area, and the right FFA. Results for the additional ANOVAs conducted separately for activity in the left HG, Broca's area and right FFA ROIs are described above.

Figure 7 demonstrates that Attention had a significant effect on brain activity in the left planum polare, the left angular gyrus, the left lingual gyrus, the right temporal pole, the right supramarginal gyrus, the right inferior parietal lobule, as well as in the ventromedial prefrontal cortex/orbitofrontal cortex and posterior cingulate bilaterally. In all these areas, activity



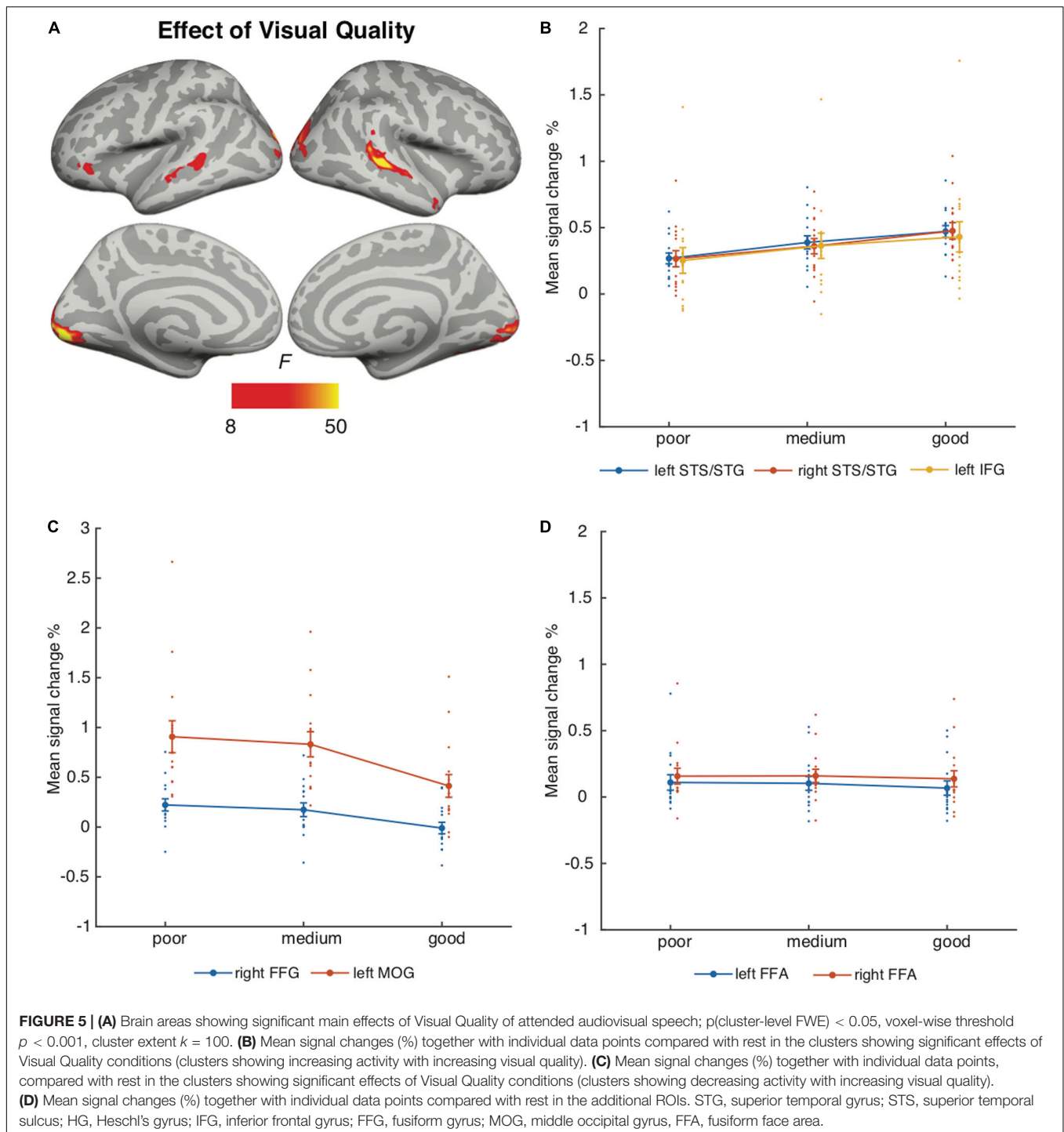
was higher during attention-to- speech than during attention-to-the-fixation cross.

For additional data, see **Supplementary Material**. **Supplementary Figure S1** presents correlations of the behavioral data with BOLD signal in selected ROIs. **Supplementary Figure S2** depicts raw BOLD activity for a sub-sample of participants in selected ROIs. **Supplementary Figure S3** depicts activity in all conditions for all participants.

Supplementary Figure S4 demonstrates individual participant data for a sub-sample of participants.

DISCUSSION

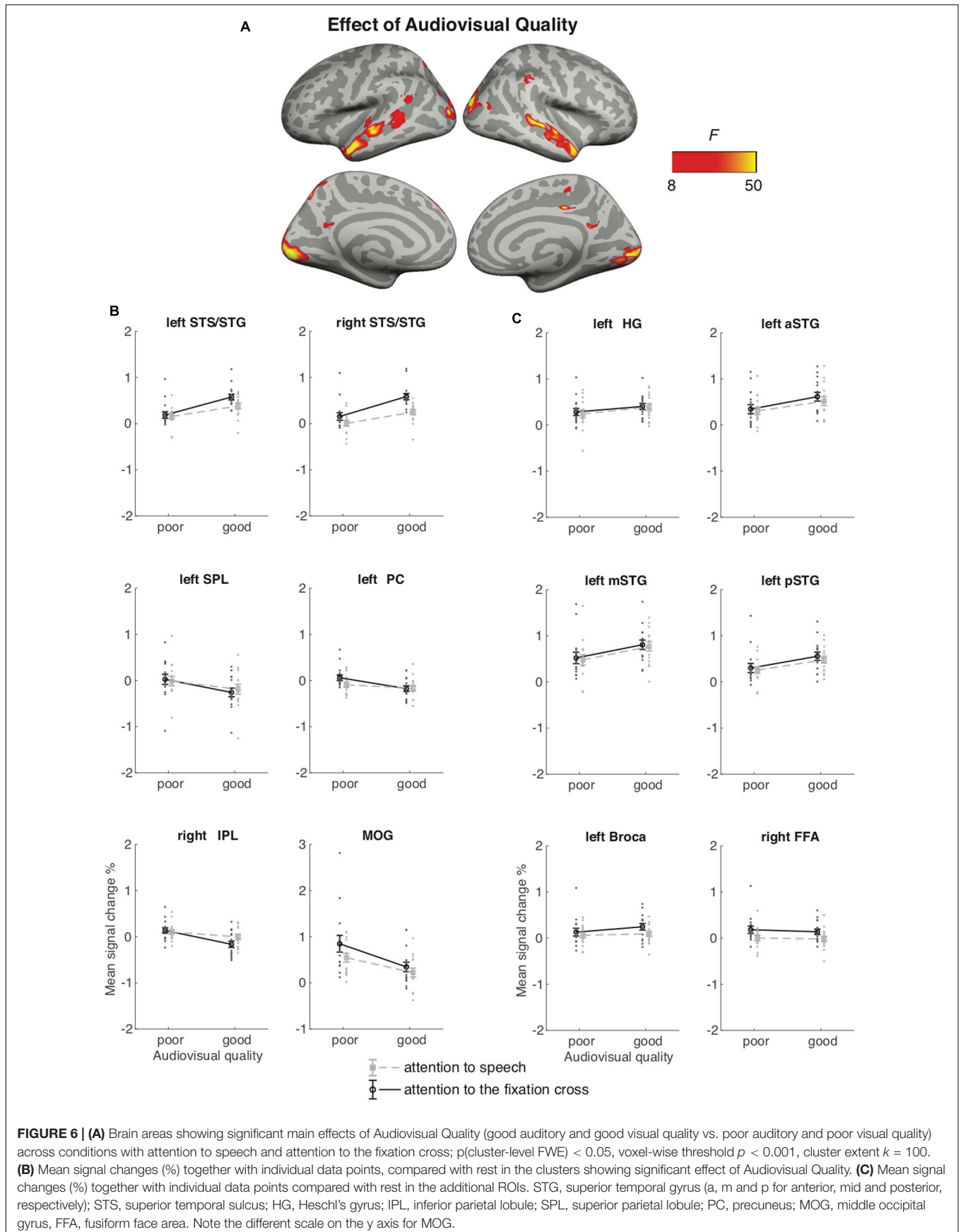
We investigated brain areas activated during selective attention to audiovisual dialogues. In particular, we expected

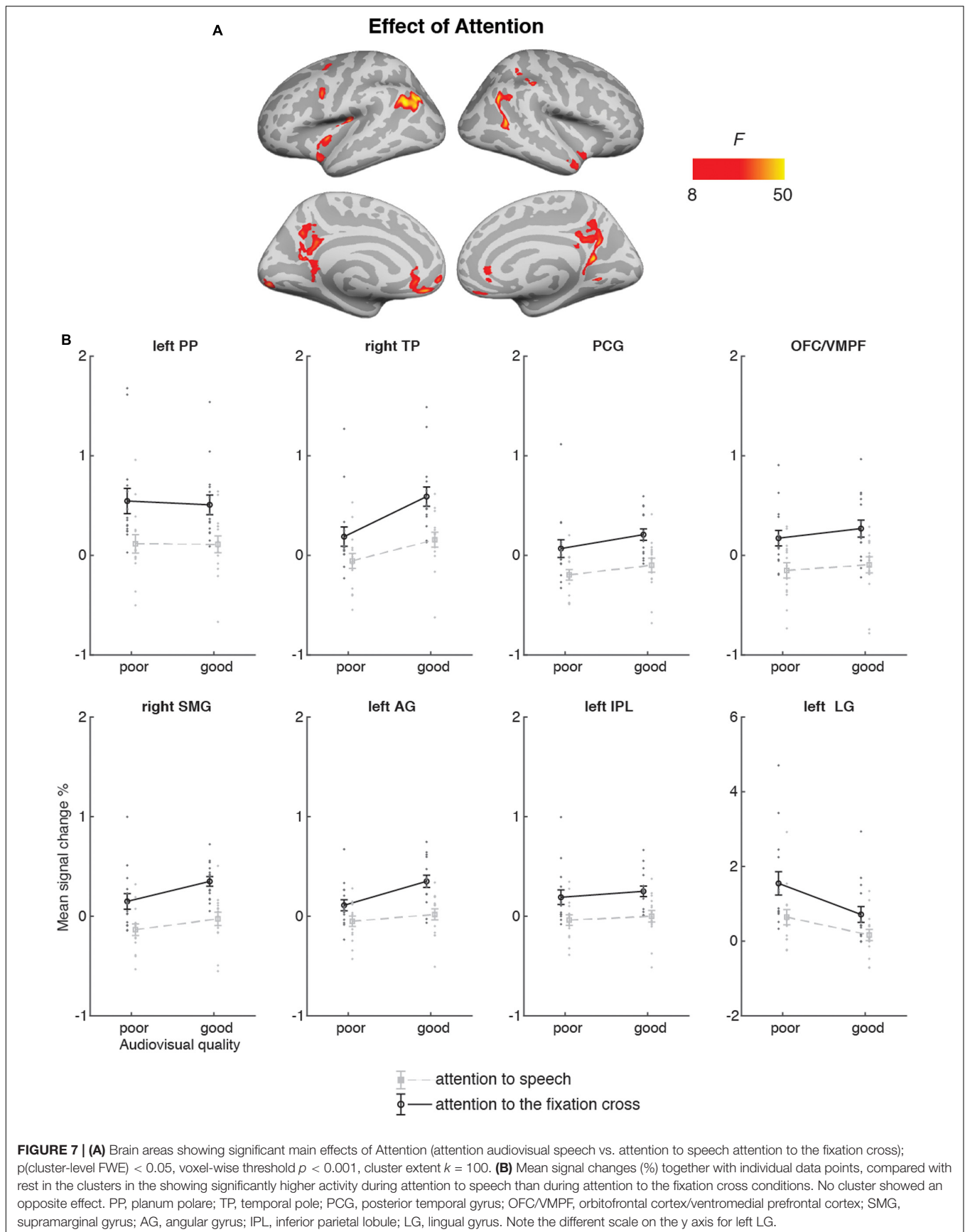


attention-related modulations in the auditory cortex and fronto-parietal activity during selective attention to naturalistic dialogues with varying auditory and visual quality. Behaviorally, we observed that increased quality of both auditory and visual information resulted in improved accuracy in answering to the questions related to the content of dialogues. Hence, expectedly, both increased auditory quality (e.g., Davis and Johnsrude, 2003) and increased visual quality (Sumby and Pollack, 1954)

facilitated speech comprehension. However, no significant interaction between Auditory Quality and Visual Quality was observed. Thus, our results are not able to give full support to maximal facilitation of speech processing by visual speech at the intermediate signal-to-noise ratio reported, for instance, by McGettigan et al. (2012) and Ross et al. (2007).

In the fMRI analysis, the main effect of Auditory Quality showed that increasing speech quality was associated with





increased activity in the (bilateral) STG/STS, which corroborates previous studies on speech intelligibility (e.g., Scott et al., 2000; Davis and Johnsrude, 2003; Obleser et al., 2007; Okada et al., 2010; McGettigan et al., 2012; Evans et al., 2014, 2016). The bilaterally enhanced activity in the STG/STS is most probably related to increased speech comprehension with increasing availability of linguistic information. The STG/STS activity extended to the temporal pole, which might be associated with enhanced semantic processing with increasing speech quality (Patterson et al., 2007). The right STG/STS activity observed here might also be related to prosodic processing during attentive listening (Alho et al., 2006; McGettigan et al., 2013; Kyong et al., 2014). The right temporal pole, in turn, its most anterior part in particular, has also been associated with social cognition (Olson et al., 2013), which may have been triggered by our naturalistic audiovisual dialogues. In addition, we observed increasing activity in the left angular gyrus and left medial frontal gyrus with increasing speech intelligibility. Enhanced activity in the left angular gyrus may reflect successful speech comprehension, stemming either from increased speech quality or from facilitated semantic processing due to improved speech quality (Humphries et al., 2007; Obleser and Kotz, 2010). The left medial frontal gyrus, in turn, has been attributed to semantic processing as a part of a semantic network (Binder et al., 2009). Hence, an increase in these activations with improving speech quality implies a successful integration of linguistic information onto the existing semantic network and improved comprehension of the spoken input – extending beyond the STG/STS.

The main effect of Visual Quality demonstrated increasing activity in the bilateral occipital cortex and right fusiform gyrus with decreasing visual quality – areas related to object and face recognition, respectively (e.g., Weiner and Zilles, 2016). Enhanced activity in these areas might be due to, for instance, noise-modulation of the videos that contained more random motion on the screen than good quality videos. Visual noise has been shown to activate primary regions in the occipital cortex more than coherent motion (e.g., Braddick et al., 2001). It is, however, also possible that viewing masked visual speech required more visual attention than viewing the unmasked visual speech, perhaps contributing to enhanced activity in the degraded visual conditions especially in the fusiform gyrus. Nevertheless, activity in the middle occipital gyrus was higher for poor visual quality combined with poor auditory quality than for good visual and auditory quality even during attention to the fixation cross (**Figure 6**). This suggests that increased visual cortex activity for poorer visual quality was at least partly caused by random motion of the masker. In a study that used blurring to decrease the reliability of visual information, stronger activity (and connectivity with the STS) was found in the extrastriate visual cortex for more reliable (i.e., less noisy) visual information (Nath and Beauchamp, 2011). It is possible that enhancements in the extrastriate visual cortex were obscured by our noise-modulations, affecting processing already in the primary visual cortex. Activity enhancements with poor (contra good) audiovisual quality were also observed in the left superior parietal lobule, precuneus and right inferior parietal

lobule in both attention conditions, implying contribution of random motion in the masker to these effects as well. Increased visual quality was also associated with enhanced activity in the bilateral STG/STS, corroborating other studies reporting these areas being involved in multisensory integration (e.g., Beauchamp et al., 2004a,b, Lee and Noppeney, 2011). Facilitation of speech processing in STG/STS areas by visual speech might be mediated by cortico-cortical connections between brain areas involved in visual and auditory speech processing (Cappe and Barone, 2005; van Wassenhove et al., 2005). We also observed an increase in the left IFG activity with increasing visual quality, an area related to the processing of high-order linguistic information (e.g., Obleser et al., 2007). Activity in the left IFG has been also associated with integration of auditory and visual speech (Lee and Noppeney, 2011), as well as speech and gestures (Willems et al., 2009), suggesting its involvement in multimodal integration also in the current study.

The 2×2 ANOVA for brain activity during attention to speech and during attention to the fixation cross showed a main effect of Audiovisual Quality in the bilateral STG/STS (extending to the temporal pole), due to higher activity for good auditory and good visual quality than for poor auditory and poor visual quality. These STG/STS effects were also observed during attention to the fixation cross, implying quite automatic bottom-up speech processing with enhanced audiovisual quality. The 2×2 ANOVA indicated also main effects of attention, that is, higher activity during attention to the dialogue than attention to the fixation cross in the left planum polare, angular and lingual gyrus, as well as the right temporal pole. We also observed activity in the dorsal part of the right inferior parietal lobule and supramarginal gyrus, as well as in the orbitofrontal/ventromedial frontal gyrus and posterior cingulate bilaterally. One might wonder why attending to the dialogues in relation to attending to the fixation cross was not associated with activity enhancements in the STG/STS as in some previous studies on selective attention to continuous speech (e.g., Alho et al., 2003, 2006). One possible explanation is the ease of the visual control task (i.e., counting the rotations of the fixation cross), eliminating the need to disregard audiovisual speech in the background altogether. This interpretation is also supported by the STG/STS activations observed even during attention to the fixation cross, at least when the audiovisual quality in the to-be-ignored dialogue was good (see **Figure 6**). Areas in the planum polare have been shown to be associated with task-related manipulations in relation to speech stimuli (Harinen et al., 2013; Wikman and Rinne, 2019).

Auditory attention effects have also been reported outside the STG/STS, for instance, in the middle and superior frontal gyri, precuneus, as well as superior parietal inferior and superior parietal lobule (e.g., Degerman et al., 2006; Salmi et al., 2007). These areas are at least partly involved in the top-down control of auditory cortex during selective attention. Interestingly, even though the participants attended to visual stimuli both during attention to speech and attention to the fixation cross, activity was higher in the lingual gyrus (approximately in areas V2/V3 of the visual cortex) during attention to speech. This effect is presumably explained by differences in visual attention between

the tasks (see, e.g., Martínez et al., 1999). In other words, while both tasks demanded visual attention, task-related processing of visual speech was presumably more attention-demanding, especially when the faces were masked, than processing of fixation-cross rotations.

It should be noted that some activity enhancements during attention-to-speech in relation to attention-to-the-fixation-cross might not be associated with attention *per se*, but with higher effort in task performance during attention to the dialogues. However, if this was the case one would expect to see an interaction between Task and Stimulus Quality for the ANOVA depicted in **Figure 6**. That is, higher effort required to process the dialogues with poor auditory quality should modulate activations during the attention-to-speech conditions but not during the attention-to-the-fixation-cross conditions, when participants did not process the dialogues. Moreover, the 3×3 ANOVA for the nine auditory-visual quality combinations during the attend-to-speech conditions (**Figures 4 and 5**) did not show enhanced activations due to poorer stimulus quality, except for the bilateral occipital cortex and right fusiform gyrus, where activity increased with decreasing quality of visual speech (see **Figure 4B**). Also, as noted above, enhanced demands for visual attention may have contributed to these effects. Alternatively, these effects might be associated with enhanced effort in perceiving visual speech with decreasing visual quality. However, as discussed above, the 2×2 ANOVA for brain activity during attention-to-speech and during attention-to-the-fixation-cross showed a main effect of Audiovisual Quality in the middle occipital gyri, with higher activity for the poor-poor than good-good auditory-visual quality combination, but no significant interaction of Audiovisual Quality and Attention. Since these activity enhancements in the middle occipital gyri and in some parietal areas, associated with decreasing audiovisual quality, were quite similar during the two attention tasks (see **Figure 6B**), it is likely that these effects were due to bottom-up effects associated with visual stimulation differences between conditions with poor and good visual quality, rather than due to differences in effort.

Unfortunately, the relatively small number of participants in the present study does not allow for investigation of behavior-brain relationships. However, it should be noted that the main findings of the present effects of attention and auditory and visual quality of speech on brain activity and performance were replicated in our two subsequent fMRI studies, which are still in preparation (for preliminary results, see Wikman et al., 2018, 2019).

In line with the previous studies on selective attention to continuous speech (Alho et al., 2003, 2006; Scott et al., 2004), attention to audiovisual dialogues did not significantly engage dorsolateral prefrontal and superior parietal areas. This may be due to high automaticity of selective listening to continuous speech, which might, hence, be quite independent of fronto-parietal attentional control (Alho et al., 2006). However, for the present audiovisual attention to speech, we observed activation in the left inferior parietal lobule, which may be related to attentive auditory processing (e.g., Rinne et al., 2009; Alain et al., 2010).

Furthermore, attention to audiovisual speech elicited enhanced activity in the orbitofrontal/ventromedial prefrontal

cortex in comparison with attention to the fixation cross. One possible explanation would be that this activity is related to processing of semantic information (e.g., Binder et al., 2009) in attended speech in contrast to visual information in the fixation cross. Alternatively, this effect may be related to the social aspect of the attended dialogues, since the ventromedial frontal area is associated with social cognition, such as theory of mind and moral judgment (Bzdok et al., 2012), as well as evaluation of other persons' traits (Araujo et al., 2013). Moreover, enhanced activity in the posterior cingulate and right superior temporal pole observed here during attention to speech may be related to social perception, as both these areas have been involved in social cognition (Bzdok et al., 2012). To our knowledge, no previous study has shown that attending to emotionally neutral dialogues would enhance activity in these three brain regions related to social perception and cognition.

To summarize, our study is the first to present findings on selective attention to natural audiovisual dialogues. Our results demonstrate that increased auditory and visual quality of speech facilitated selective listening to the dialogues, seen in enhanced brain activity in the bilateral STG/STS and the temporal pole. Enhanced activity in the temporal pole might be related to semantic processing particularly in the left hemisphere, whereas in the right hemisphere, it may index processing of social information activated during attention to the dialogues. The fronto-parietal network was associated with enhanced activity during attention to speech, reflecting top-down attentional control. Attention to audiovisual speech also activated the orbitofrontal/ventromedial prefrontal cortex – a region associated with social and semantic cognition. Hence, our findings on selective attention in realistic audiovisual dialogues emphasize not only involvement of brain networks related to audiovisual speech processing and semantic comprehension but, as a novel observation, the social brain network.

DATA AVAILABILITY STATEMENT

The datasets generated for this study are available on request to the corresponding author.

ETHICS STATEMENT

The studies involving human participants were reviewed and approved by Ethics Review Board in the Humanities and Social and Behavioral Sciences, University of Helsinki. The patients/participants provided their written informed consent to participate in this study. Written informed consent was obtained from the individual(s) for the publication of any potentially identifiable images or data included in this manuscript.

AUTHOR CONTRIBUTIONS

AL, KA, MV, and VS designed the fMRI experiment. AL, MV, and VS prepared the stimuli. MV collected the fMRI data.

AL and MV performed the analysis. AL, MV, and KA wrote the manuscript in collaboration with PW, VS, and MM. MM, VS, and PW contributed to the fMRI data analysis.

FUNDING

This study was supported by the Academy of Finland (Grant No. 297848) and by the Russian Science Foundation (grant code RSF 19-18-00534). MV was awarded with Erasmus + scholarship.

ACKNOWLEDGMENTS

We thank Artturi Ylinen for help with stimulus preparation and data collection and Miika Leminen for methodological support.

REFERENCES

- Alain, C., Shen, D., Yu, H., and Grady, C. (2010). Dissociable memory- and response-related activity in parietal cortex during auditory spatial working memory. *Front. Psychol.* 1:202. doi: 10.3389/fpsyg.2010.00202
- Alho, K., Medvedev, S. V., Pakhomov, S. V., Roudas, M. S., Tervaniemi, M., Reinikainen, K., et al. (1999). Selective tuning of the left and right auditory cortices during spatially directed attention. *Cogn. Brain Res.* 7, 335–341. doi: 10.1016/s0926-6410(98)00036-6
- Alho, K., Rinne, T., Herron, T. J., and Woods, D. L. (2014). Stimulus-dependent activations and attention-related modulations in the auditory cortex: a meta-analysis of fMRI studies. *Hear. Res.* 307, 29–41. doi: 10.1016/j.heares.2013.08.001
- Alho, K., Salmi, J., Koistinen, S., Salonen, O., and Rinne, T. (2015). Top-down controlled and bottom-up triggered orienting of auditory attention to pitch activate overlapping brain networks. *Brain Res.* 1626, 136–145. doi: 10.1016/j.brainres.2014.12.050
- Alho, K., Vorobyev, V. A., Medvedev, S. V., Pakhomov, S. V., Roudas, M. S., Tervaniemi, M., et al. (2003). Hemispheric lateralization of cerebral blood-flow changes during selective listening to dichotically presented continuous speech. *Cogn. Brain Res.* 17, 201–211. doi: 10.1016/s0926-6410(03)00091-0
- Alho, K., Vorobyev, V. A., Medvedev, S. V., Pakhomov, S. V., Starchenko, M. G., Tervaniemi, M., et al. (2006). Selective attention to human voice enhances brain activity bilaterally in the superior temporal sulcus. *Brain Res.* 1075, 142–150. doi: 10.1016/j.brainres.2005.11.103
- Araujo, H. F., Kaplan, J., and Damasio, A. (2013). Cortical midline structures and autobiographical-self processes: an activation-likelihood estimation meta-analysis. *Front. Hum. Neurosci.* 7:548. doi: 10.3389/fnhum.2013.00548
- Beauchamp, M. S., Argall, B. D., Bodurka, J., Duyn, J. H., and Martin, A. (2004a). Unraveling multisensory integration: patchy organization within human STS multisensory cortex. *Nat. Neurosci.* 7, 1190–1192. doi: 10.1038/nn1333
- Beauchamp, M. S., Lee, K. E., Argall, B. D., and Martin, A. (2004b). Integration of auditory and visual information about objects in superior temporal sulcus. *Neuron* 41, 809–823. doi: 10.1016/s0896-6273(04)00070-4
- Binder, J. R., Desai, R. H., Graves, W. W., and Conant, L. L. (2009). Where is the semantic system? A critical review and meta-analysis of 120 functional neuroimaging studies. *Cereb. Cortex* 19, 2767–2796. doi: 10.1093/cercor/bhp055
- Bishop, C. W., and Miller, L. M. (2009). A multisensory cortical network for understanding speech in noise. *J. Cogn. Neurosci.* 21, 1790–1805. doi: 10.1162/jocn.2009.21118
- Boersma, P., and Weenink, D. (2001). Praat: doing phonetics by computer. *Glott Int.* 5, 341–345. doi: 10.1016/j.jvoice.2019.07.004
- Braddick, O. J., O'Brien, J. M., Wattam-Bell, J., Atkinson, J., Hartley, T., and Turner, R. (2001). Brain areas sensitive to coherent visual motion. *Perception* 30, 61–72. doi: 10.1068/p3048

SUPPLEMENTARY MATERIAL

The Supplementary Material for this article can be found online at: <https://www.frontiersin.org/articles/10.3389/fnins.2020.00436/full#supplementary-material>

FIGURE S1 | Correlations of the behavioral data with BOLD signal in selected ROIs. Y-axis: behavioral quiz scores (1–7), x-axis: BOLD signal change (%). A1–A3 denote poor, medium and good auditory quality; V1–V3 denote poor, medium and good visual quality. Each experimental condition is depicted in different color. STG, superior temporal gyrus (a, m, and p for anterior, mid and posterior, respectively); STS, superior temporal sulcus; HG, Heschl's gyrus; FFA, fusiform face area.

FIGURE S2 | Raw BOLD activity (first run) for five participants in selected ROIs. Y-axis: time, volume, x-axis: mean BOLD, arbitrary unit. STG, superior temporal gyrus (a, m, and p for anterior, mid and posterior, respectively); STS, superior temporal sulcus; HG, Heschl's gyrus; FFA, fusiform face area.

FIGURE S3 | The average signal changes across all conditions.

FIGURE S4 | Individual participant data for 5 participants across all conditions.

- Bzdok, D., Schilbach, L., Vogeley, K., Schneider, K., Laird, A. R., Langner, R., et al. (2012). Parsing the neural correlates of moral cognition: ALE meta-analysis on morality, theory of mind, and empathy. *Brain Struct. Funct.* 217, 783–796. doi: 10.1007/s00429-012-0380-y
- Cappe, C., and Barone, B. (2005). Heteromodal connections supporting multisensory integration at low levels of cortical processing in the monkey. *Eur. J. Neurosci.* 22, 2886–2902. doi: 10.1111/j.1460-9568.2005.04462.x
- Davis, M. H., and Johnsrude, I. S. (2003). Hierarchical processing in spoken language comprehension. *J. Neurosci.* 23, 3423–3431. doi: 10.1523/JNEUROSCI.23-08-03423.2003
- Degerman, A., Rinne, T., Salmi, J., Salonen, O., and Alho, K. (2006). Selective attention to sound location or pitch studied with fMRI. *Brain Res.* 1077, 123–134. doi: 10.1016/j.brainres.2006.01.025
- Evans, S., Kyong, J. S., Rosen, S., Golestani, N., Warren, J. E., McGettigan, C., et al. (2014). The pathways for intelligible speech: multivariate and univariate perspectives. *Cereb. Cortex* 24, 2350–2361. doi: 10.1093/cercor/bht083
- Evans, S., McGettigan, C., Agnew, Z. K., Rosen, S., and Scott, S. K. (2016). Getting the cocktail party started: masking effects in speech perception. *J. Cogn. Neurosci.* 28, 483–500. doi: 10.1162/jocn_a_00913
- Greenwood, D. D. (1990). A cochlear frequency-position function for several species—29 years later. *J. Acoust. Soc. Am.* 87, 2592–2605. doi: 10.1121/1.399052
- Grill-Spector, K., and Weiner, K. S. (2014). The functional architecture of the ventral temporal cortex and its role in categorization. *Nat. Rev. Neurosci.* 15, 536–548. doi: 10.1038/nrn3747
- Harinen, K., Aaltonen, O., Salo, E., Salonen, O., and Rinne, T. (2013). Task-dependent activations of human auditory cortex to prototypical and nonprototypical vowels. *Hum. Brain Mapp.* 34, 1272–1281. doi: 10.1002/hbm.21506
- Hill, K. T., and Miller, L. M. (2010). Auditory attentional control and selection during cocktail party listening. *Cereb. Cortex* 20, 583–590. doi: 10.1093/cercor/bhp124
- Humphries, C., Binder, J. R., Medler, D. A., and Liebenthal, E. (2007). Time course of semantic processes during sentence comprehension: an fMRI study. *NeuroImage* 36, 924–932. doi: 10.1016/j.neuroimage.2007.03.059
- Kyong, J. S., Scott, S. K., Rosen, S., Howe, T. B., Agnew, Z. K., and McGettigan, C. (2014). Exploring the roles of spectral detail and intonation contour in speech intelligibility: an fMRI study. *J. Cogn. Neurosci.* 26, 1748–1763. doi: 10.1162/jocn_a_00583
- Lee, H., and Noppeney, U. (2011). Physical and perceptual factors shape the neural mechanisms that integrate audiovisual signals in speech comprehension. *J. Neurosci.* 31, 11338–11350. doi: 10.1523/JNEUROSCI.6510-10.2011
- Liakakis, G., Nickel, J., and Seitz, R. J. (2011). Diversity of the inferior frontal gyrus—a meta-analysis of neuroimaging studies. *Behav. Brain Res.* 225, 341–347. doi: 10.1016/j.bbr.2011.06.022
- Liebenthal, E., Desai, R. H., Humphries, C., Sabri, M., and Desai, A. (2014). The functional organization of the left STS: a large scale meta-analysis of PET and

- fMRI studies of healthy adults. *Front. Neurosci.* 8:289. doi: 10.3389/fnins.2014.00289
- Martinez, A., Anllo-Vento, L., Sereno, M. I., Frank, L. R., Buxton, R. B., Dubowitz, D. J., et al. (1999). Involvement of striate and extrastriate visual cortical areas in spatial attention. *Nat. Neurosci.* 2, 364–369. doi: 10.1038/7274
- McFarquhar, M., McKie, S., Emsley, R., Suckling, J., Elliott, R., and Williams, S. (2016). Multivariate and repeated measures (MRM): a new toolbox for dependent and multimodal group-level neuroimaging data. *NeuroImage* 132, 373–389. doi: 10.1016/j.neuroimage.2016.02.053
- McGettigan, C., Eisner, F., Agnew, Z. K., Manly, T., Wisbey, D., and Scott, S. K. (2013). T'ain't what you say, it's the way that you say it – Left insula and inferior frontal cortex work in interaction with superior temporal regions to control the performance of vocal impersonations. *J. Cogn. Neurosci.* 25, 1875–1886. doi: 10.1162/jocn_a_00427
- McGettigan, C., Faulkner, A., Altarelli, I., Obleser, J., Baverstock, H., and Scott, S. K. (2012). Speech comprehension aided by multiple modalities: behavioural and neural interactions. *Neuropsychologia* 50, 762–776. doi: 10.1016/j.neuropsychologia.2012.01.010
- Mesgarani, N., and Chang, E. F. (2012). Selective cortical representation of attended speaker in multi-talker speech perception. *Nature* 485, 233–236. doi: 10.1038/nature11020
- Nath, A. R., and Beauchamp, M. S. (2011). Dynamic changes in superior temporal sulcus connectivity during perception of noisy audiovisual speech. *J. Neurosci.* 31, 1704–1714. doi: 10.1523/JNEUROSCI.4853-10.2011
- Obleser, J., and Kotz, S. A. (2010). Expectancy constraints in degraded speech modulate the language comprehension network. *Cereb. Cortex* 20, 633–640. doi: 10.1093/cercor/bhp128
- Obleser, J., Wise, R. J., Dresner, M. A., and Scott, S. K. (2007). Functional integration across brain regions improves speech perception under adverse listening conditions. *J. Neurosci.* 27, 2283–2289. doi: 10.1523/JNEUROSCI.4663-06.2007
- Okada, K., Rong, F., Venezia, J., Matchin, W., Hsieh, I. H., Saberi, K., et al. (2010). Hierarchical organization of human auditory cortex: evidence from acoustic invariance in the response to intelligible speech. *Cereb. Cortex* 20, 2486–2495. doi: 10.1093/cercor/bhp318
- Oldfield, R. C. (1971). The assessment and analysis of handedness: the Edinburgh inventory. *Neuropsychologia* 9, 97–113. doi: 10.1016/0028-3932(71)90067-4
- Olson, I. R., McCoy, D., Klobusicky, E., and Ross, L. A. (2013). Social cognition and the anterior temporal lobes: a review and theoretical framework. *Soc. Cogn. Affect. Neurosci.* 8, 123–133. doi: 10.1093/scan/nss119
- O'Sullivan, J. A., Power, A. J., Mesgarani, N., Rajaram, S., Foxe, J. J., Shinn-Cunningham, B. G., et al. (2015). Attentional selection in a cocktail party environment can be decoded from single-trial EEG. *Cereb. Cortex* 25, 1697–1706. doi: 10.1093/cercor/bht355
- Patterson, K., Nestor, P. J., and Rogers, T. T. (2007). Where do you know what you know? The representation of semantic knowledge in the human brain. *Nat. Rev. Neurosci.* 8, 976–987. doi: 10.1038/nrn2277
- Puschmann, S., Steinkamp, S., Gillich, I., Mirkovic, B., Debener, S., and Thiel, C. M. (2017). The right temporoparietal junction supports speech tracking during selective listening: evidence from concurrent EEG-fMRI. *J. Neurosci.* 37, 11505–11516. doi: 10.1523/JNEUROSCI.1007-17.2017
- Rimmele, J. M., Zion Golumbic, E., Schroger, E., and Poeppel, D. (2015). The effects of selective attention and speech acoustics on neural speech-tracking in a multi-talker scene. *Cortex* 68, 144–154. doi: 10.1016/j.cortex.2014.12.014
- Rinne, T., Koistinen, S., Salonen, O., and Alho, K. (2009). Task-dependent activations of human auditory cortex during pitch discrimination and pitch memory tasks. *J. Neurosci.* 29, 13338–13343. doi: 10.1523/JNEUROSCI.3012-09.2009
- Ross, L. A., Saint-Amour, D., Leavitt, V. M., Javitt, D. C., and Foxe, J. J. (2007). Do you see what I am saying? Exploring visual enhancement of speech comprehension in noisy environments. *Cereb. Cortex* 17, 1147–1153. doi: 10.1093/cercor/bhl024
- Salmi, J., Rinne, T., Degerman, A., Salonen, O., and Alho, K. (2007). Orienting and maintenance of spatial attention in audition and vision: multimodal and modality-specific brain activations. *Brain Struct. Funct.* 212, 181–194. doi: 10.1007/s00429-007-0152-2
- Scott, S. K., Blank, C. C., Rosen, S., and Wise, R. J. (2000). Identification of a pathway for intelligible speech in the left temporal lobe. *Brain* 123, 2400–2406. doi: 10.1093/brain/123.12.2400
- Scott, S. K., Rosen, S., Wickham, L., and Wise, R. J. (2004). A positron emission tomography study of the neural basis of informational and energetic masking effects in speech perception. *J. Acoust. Soc. Am.* 115, 813–821. doi: 10.1121/1.1639336
- Shannon, R. V., Zeng, F. G., Kamath, V., Wygonski, J., and Ekelid, M. (1995). Speech recognition with primarily temporal cues. *Science* 270, 303–304. doi: 10.1126/science.270.5234.303
- Sumbly, W. H., and Pollack, I. (1954). Visual contribution to speech intelligibility in noise. *J. Acoust. Soc. Am.* 25, 212–215.
- Tzourio, N., Massiou, F. E., Crivello, F., Joliot, M., Renault, B., and Mazoyer, B. (1997). Functional anatomy of human auditory attention studied with PET. *NeuroImage* 5, 63–77. doi: 10.1006/nimg.1996.0252
- van Wassenhove, V., Grant, K. W., and Poeppel, D. (2005). Visual speech speeds up the neural processing of auditory speech. *Proc. Nat. Acad. Sci. U.S.A.* 102, 41181–41186. doi: 10.1073/pnas.0408949102
- Weiner, K. S., and Zilles, K. (2016). The anatomical and functional specialization of the fusiform gyrus. *Neuropsychologia* 83, 48–62. doi: 10.1016/j.neuropsychologia.2015.06.033
- Whitfield-Gabrieli, S., and Nieto-Castanon, A. (2012). Conn: a functional connectivity toolbox for correlated and anticorrelated brain networks. *Brain Connect.* 2, 125–141. doi: 10.1089/brain.2012.0073
- Wikman, P., and Rinne, T. (2019). Interaction of the effects associated with auditory-motor integration and attention-engaging listening tasks. *Neuropsychologia* 124, 322–336. doi: 10.1016/j.neuropsychologia.2018.11.006
- Wikman, P., Sahari, E., Leminen, A., Leminen, M., Laine, M., and Alho, K. (2018). “Selective auditory attention during naturalistic audio-visual dialogues with shuffled lines,” in *Neurobiology of Speech and Language, Proceedings of the 2nd International Workshop*, eds O. Shcherbakova and Y. Shtyrov St. Petersburg, 60.
- Wikman, P., Ylinen, A. Y. H., Leminen, A., Leminen, M., and Alho, K. (2019). “Task-related effects during natural audio-visual dialogues in fMRI,” in : *Proceedings of the 11th Annual Meeting of the Society for the Neurobiology of Language*, Helsinki, 20–22.
- Wild, C. J., Yusuf, A., Wilson, D. E., Peelle, J. E., Davis, M. H., and Johnsrude, I. S. (2012). Effortful listening: the processing of degraded speech depends critically on attention. *J. Neurosci.* 32, 14010–14021. doi: 10.1523/JNEUROSCI.1528-12.2012
- Willems, R. M., Ozyurek, A., and Hagoort, P. (2009). Differential roles for left inferior frontal and superior temporal cortex in multimodal integration of action and language. *NeuroImage* 47, 1992–2004. doi: 10.1016/j.neuroimage.2009.05.066
- Zatorre, R. J., Mondor, T. A., and Evans, A. C. (1999). Auditory attention to space and frequency activates similar cerebral systems. *NeuroImage* 10, 544–554. doi: 10.1006/nimg.1999.0491
- Zion Golumbic, E. M., Ding, N., Bickel, S., Lakatos, P., Schevon, C. A., McKhann, G. M., et al. (2013). Mechanisms underlying selective neuronal tracking of attended speech at a “cocktail party”. *Neuron* 77, 980–991. doi: 10.1016/j.neuron.2012.12.037

Conflict of Interest: The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

Copyright © 2020 Leminen, Verwoert, Moiala, Salmela, Wikman and Alho. This is an open-access article distributed under the terms of the Creative Commons Attribution License (CC BY). The use, distribution or reproduction in other forums is permitted, provided the original author(s) and the copyright owner(s) are credited and that the original publication in this journal is cited, in accordance with accepted academic practice. No use, distribution or reproduction is permitted which does not comply with these terms.