

AVbook, a high-frame-rate corpus of narrative audiovisual speech for investigating multimodal speech perception

Enrico Varano¹, Tobias Reichenbach^{1,2,*}

1 Department of Bioengineering and Centre for Neurotechnology, Imperial College London, London, United Kingdom.

2 Department of Artificial Intelligence in Biomedical Engineering, Friedrich-Alexander-University Erlangen-Nuremberg, Erlangen, Germany. tobias.j.reichenbach@fau.de

* Contact author: tobias.j.reichenbach@fau.de

Abstract

We present an audiovisual speech corpus that is designed for cognitive neuroscience studies and that can also be employed for research on audiovisual speech recognition. The corpus consists of 3.6 hours of audiovisual recordings of two speakers, one male and one female, reading passages from a narrative English text. The visual recordings were acquired at a high frame rate of 119.88 frames per second (fps) and exported at a high resolution of 528 x 718 pixels. The speech is pronounced with a neutral British accent and is directed at the camera. Both speakers read the same 59 passages of a book, for a total of 1h50' each. The passage scripts, largely contiguous within a non-fiction source book chosen for its compelling content, were selected and lightly edited to keep subjects who might listen to it interested and alert. As tools to test comprehension and attention, sets of four multiple-choice questions were written for each passage. A short written summary is also provided for each recording. To enable audiovisual synchronisation when presenting the stimuli, four videos of an electronic clapperboard were recorded in line with the corpus. Stimulus synchronisation of 0 ± 4 ms was achieved by pairing these with a high frame rate commercial monitor and a photo-sensor. The audiovisual speech material, the corresponding text, synchronization material, comprehension questions and written summaries set are available on the web for research use.

Introduction

Visual cues play a prominent role in natural communication and have been shown to improve speech comprehension with and without background noise (Reisberg et al., 1987, Ross et al., 2007). This audiovisual benefit occurs across the scale of linguistic units, from syllables (Bernstein et al., 2004) to words (Sumbly and Pollack, 1954) and sentences (Grant and Seitz, 2000), and extends to automatic speech recognisers (Matthews et al., 2002).

Databases of audiovisual (AV) speech material play an important role in enabling studies on the neural mechanisms that underlie audiovisual integration regarding speech perception. Studies on this issue traditionally employed brain imaging techniques to investigate participants' responses to short speech tokens such as phonemes and syllables. Neural responses to such short stimuli were commonly analysed by averaging hundreds of response trials time-aligned to the onset of their repetitive stimuli to obtain event-related potentials (ERPs; Luck, 2004). The need

for stimuli suitable for such paradigms could be met by the researchers recording the material themselves (McGurk and MacDonald, 1976; Brown et al., 2018; Sumbly and Pollack, 1954) or by employing one of many corpora published to support automatic speech recognition efforts.

The first audiovisual speech datasets, including TULIPS1 (Movellan, 1995) and AVletters (Matthews et al., 1998), accordingly consisted of few speakers reading digits and letters. Later datasets, for instance XM2VTS (Messer et al., 1999), AVICAR (Lee et al., 2004), VidTIMIT (Sanderson and Lovell, 2009), GRID (Cook et al., 2006), BL (Benezeth and Bachman, 2011) or MODALITY (Czyzewskiet al., 2017), provided complete sentences, a larger volume of recordings, and included several additional features such as high frame rates, head poses or lip highlighting. These databases were primarily designed to support the development of algorithms for audiovisual speech processing and recognition, speaker identification and detection, affective state recognition or talking head generation.

These audiovisual speech corpora have been used extensively to study neural mechanisms of audiovisual speech processing. However, the continuous and complex nature of natural speech is not entirely reflected in the short stimuli and limited lexicon of these speech materials (Sonkusare et al., 2019). Recent advances in analysis methodologies and computational power have allowed researchers to investigate neural responses to increasingly complex stimuli including ongoing natural speech (Marmarelis, 2004; Ding and Simon, 2014; Ding and Simon, 2012; Crosse et al., 2016). These studies have typically employed audiobooks as a practical solution to present ecologically valid speech to participants.

Continuous audiovisual speech corpora are, however, less common due to the resources required to assemble and process the material (Chitu and Rothkrantz, 2012). The most commonly employed option to date is a series of weekly addresses made by a well-known male talker speaking on contemporary social, political and economic issues (Supasorn et al., 2017; O'Sullivan et al., 2017). While several hours of audiovisual speech are available, the framing is not consistent and the content, tone and familiarity of the speaker could impact the obtained results in unintended ways.

Furthermore, when presenting multimodal material to subjects in an EEG or MEG study, the precise synchronisation of the different sensory streams must be tightly controlled to take advantage of the high temporal resolution of the recordings (Schultz et al., 2020; Crosse et al., 2014). The low frame rate of 30 frames per second (fps) of the majority of corpora, including the aforementioned male speaker corpus but with the notable exception of the short sentence MODALITY corpus, limits the precision in audiovisual alignment to 33 ms.

The AVbook corpus that we present here is a narrative audiovisual speech corpus that significantly extends the limited current pool of continuous speech corpora while providing additional features aimed at facilitating speech and neuroscience research. The material consists of 3.6 hours of high frame rate and high resolution recordings of two trained speakers, one male and one female, reading the same, curated narrative text. The framing is consistent due to the employment of a teleprompter and the speakers voice is recorded with a professional clip-on microphone, yielding high quality recordings that are also suitable for audio-only experiments. This may, in turn, enable closer comparison between uni-modal and multi-modal speech processing.

Corpus

Content Design

With the experimental participant's enjoyment and attention in mind, Alfred Lansing's book "*Endurance: Shackleton's Incredible Voyage to the Antarctic*" was selected as the source for the corpus content for its compelling narrative. The first eight chapters of the book, constituting

Part 1, were divided into 59 passages with an average word count of 335. The original text was lightly edited to partition it into contextually complete passages and to remove outdated vocabulary and direct speech. The latter action was taken to aid the speakers in producing a neutral tone without overly salient parts.

Collection

The audiovisual recordings were obtained in a studio at Imperial College London, UK, using a Sony A7S II camera (Sony Corporation, Japan) mounted in a teleprompter. The audio material was recorded concurrently through the camera's auxiliary sound port connected to a Sennheiser EW100 clip-on radio microphone (Sennheiser, Germany). The speakers wore the microphone on their clothing and sat on a chair at a distance of three meters from the teleprompter and the camera. The camera was oriented in landscape and framed the speaker's head, shoulders, and chest against a grey background.

The camera was set to record at the frame rate of 119.88 fps and with a resolution of 1280 x 720 pixels. The single channel audio was recorded at 48 kHz with 32-bit resolution.

The speakers, both professional actors, one female and one male, were recruited through the website StarNow.com and selected based on a video audition for their neutral accent and clear speech. They were directed to keep their head still facing the teleprompter, and to speak with a neutral tone of voice. The scroll speed of the teleprompter was adjusted to suit the speaker's preferred pace. Mistakes during production were re-recorded following a pause to be later cut in post-processing. Taking the post-production process of removing pauses and mistakes into account, this resulted in an average passage length of 111 and 104 seconds, and an average speech rate of 182 and 192 words per minute, for the female and male speakers respectively.

Post-processing

Editing of the footage was performed on Adobe Premiere CC 2019 (Adobe Inc., USA) with the aim of cutting speech production and pronunciation mistakes, and unnatural hesitations. Pauses over 0.8 seconds were cut down to the same threshold. Video jumps and audio clicks due to cuts were smoothed out with filters and transitions, and care was taken to minimise the overlap of these transitions with speech production. The framing was cropped around the speaker's face and neck to a vertical orientation with a resolution of 528 x 718 pixels. Five frames from each speaker are shown in Figure 1.

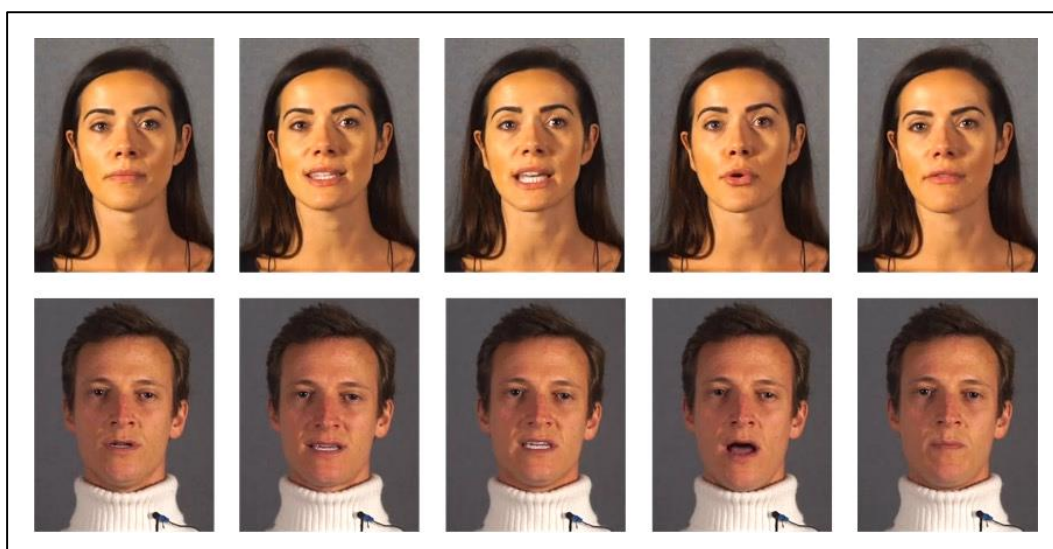


Figure 1 — Example frames from the AVbook corpus. The upper row shows five frames taken from the female speaker, and the bottom row five frames from the male speaker.

The resulting 59 videos were exported into mp4 container files using the h.265 video codec and the AAC audio codec at a frame rate of approximately 119.88 fps and 48 kHz with 16-bit resolution, respectively. The precise frame rate of the video recording was 120/1.001 fps, with each frame lasting about 8.34 ms. No clipping was observed in the audio and no further processing was performed.

The audio was also separately exported as a WAVE file. An a priori signal-to-noise ratio (SNR) was estimated using a voice activity detection (VAD) script (Brookes, 2022) following the ITU-T P.56 standard (ITU-T, 2011). The average SNR was found to be 34.0 ± 2.9 dB and 44.5 ± 7.1 dB for the female and male speakers respectively (mean and standard deviation). Similar results were obtained when employing a custom-tuned VAD script or by calculating the SNR by comparing silent video segments which were cut out of the final corpus (less than 1 dB difference in both cases for both speakers).

Script, comprehension questions and summaries

The precise wording pronounced by each speaker was checked against the script by hand. A few inconsistencies and mispronunciations that could not be corrected during the editing process were annotated in the teleprompter script file.

A short summary text was also produced for each video segment. This text gave an overview of the content of the respective passage.

Furthermore, a set of multiple-choice questions and answers was written for each passage. Four questions, each with one correct and two incorrect options, were carefully drafted to avoid revealing the answers to other questions while minimising the number of correct answers identifiable by general knowledge, context or information contained in previous passages.

Synchronisation

Synchronisation material

To enable the precise synchronisation of the audio and video streams when presenting the audiovisual material to subjects, four videos of an electronic clapperboard were recorded in the same conditions and with the same equipment as the audiovisual speech material. These 60 second synchronisation videos were post-processed and exported in the same manner as the audiovisual speech material.

The electronic clapperboard, depicted schematically in Figure 2, consisted of a standard red LED and a passive breadboard buzzer, both powered in parallel by a GW Instek GFG-8219A signal generator producing a 2.5 kHz sine wave. A mechanical switch connected in series with the signal generator was operated irregularly to generate a non-periodic on/off signal to power the diode and buzzer intermittently. The buzzer was found to have a response time of 0.3 ms and LEDs are known to have response times of the order of single-digit nanoseconds. The image of the LED and the recording of the buzzer sound contained in the synchronisation videos can therefore serve to align the audio and the video signals of the AVbook corpus.

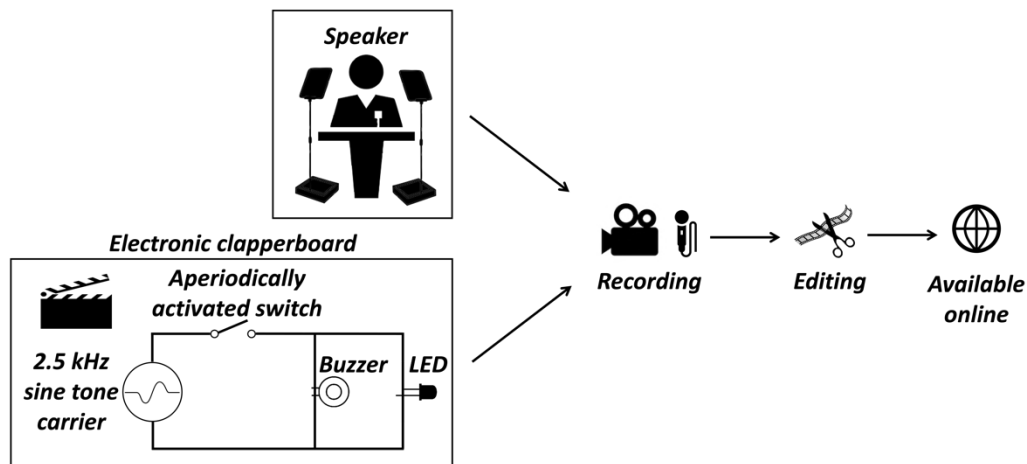


Figure 2 — Schematic depiction of the methodology for producing the audiovisual speech corpus (top) as well as the synchronisation material obtained through an electronic clapperboard (bottom).

Audiovisual alignment during presentation of the AVbook

When employing the AVbook material in behavioural or neuroscientific studies, the audiovisual synchronisation of the stimulus is essential. A diagram of the proposed stimuli presentation and alignment solution that we tested is shown in Figure 3. A custom-written stimulus presentation software module was employed to command the playback of the synchronisation files. The brightness of the image of the red LED was recorded simultaneously to the presented audio stream containing the buzzer sound and a cross-correlation was performed to determine the latency between the visual and the audio stimulus. A corresponding opposite latency shift can then be applied in the presentation software so that the resulting visual and audio stimuli are aligned.

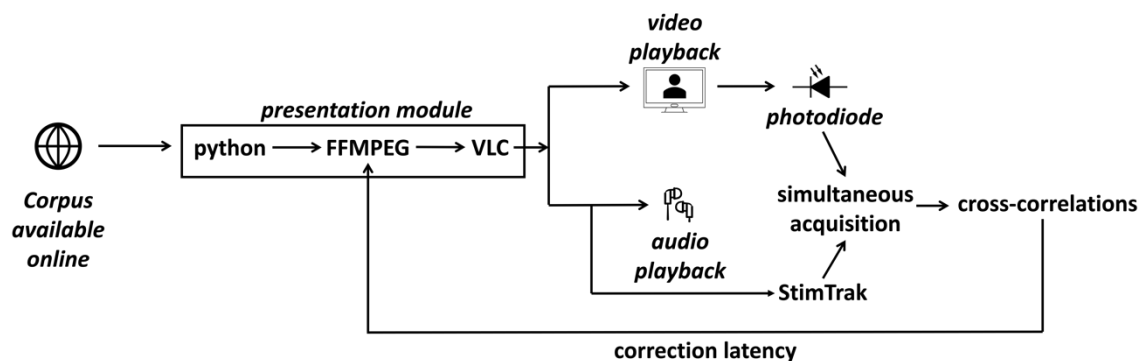


Figure 3 — Schematic depiction of the equipment and pipeline for presenting the audiovisual stimuli and for measuring and the latency between the visual and the audio signal.

Presentation hardware and module

In a test of the presentation setup, the audio stimulus was delivered diotically at a level of 70 dB(A) SPL using ER-3C insert earphones (Etymotic, USA) through a high-performance sound card (Xonar Essence STX, Asus, USA). The sound level was calibrated with a Type 4157 ear simulator (Brüel & Kjær, Denmark). The video component was delivered via a 144 Hz, 24-inch flat-screen monitor (24GM79G, LG, South Korea) set at a refresh rate of 119.88 Hz.

During an initial piloting phase, a Python 3.7 script calling the *python-vlc* library¹ was found to be the most reliable way to decode and present videos in the h.265 video codec synchronously to the audio material. Any residual timing difference between the audio and the video stimulus can then be corrected using the *-itoffset* FFmpeg flag² called through the *subprocess* python module.

Recording hardware

An actiCHamp electrophysiological recording amplifier (BrainProducts, Germany) was used to record the image of the red LED and the sound of the buzzer in the synchronisation videos. A photodiode (Photo Sensor, BrainProducts, Germany) and an acoustic adaptor (StimTrak, BrainProducts, Germany) were plugged into the auxiliary ports of the integrated amplifier. The combined data stream was then acquired through PyCorder (BrainProducts, Germany) at a sampling rate of 10 kHz, therefore allowing for temporal alignment of the audiovisual signal with a precision of 0.1 ms.

Results

A sliding window cross-correlation analysis (Figure 4) of the photodiode and acoustic adaptor signals for the synchronisation videos was performed to determine the delay of the audio signal with respect to the video stimulus over the runtime of the video. It was found that the delay averaged 41.7 ms, was stable within 0.8 ms across time and jittered by a maximum of 8.1 ms between different trials. Correcting for the mean delay in the presentation software (calling the *itoffset* ffmpeg flag) resulted in a mean delay of 0.6 ms.

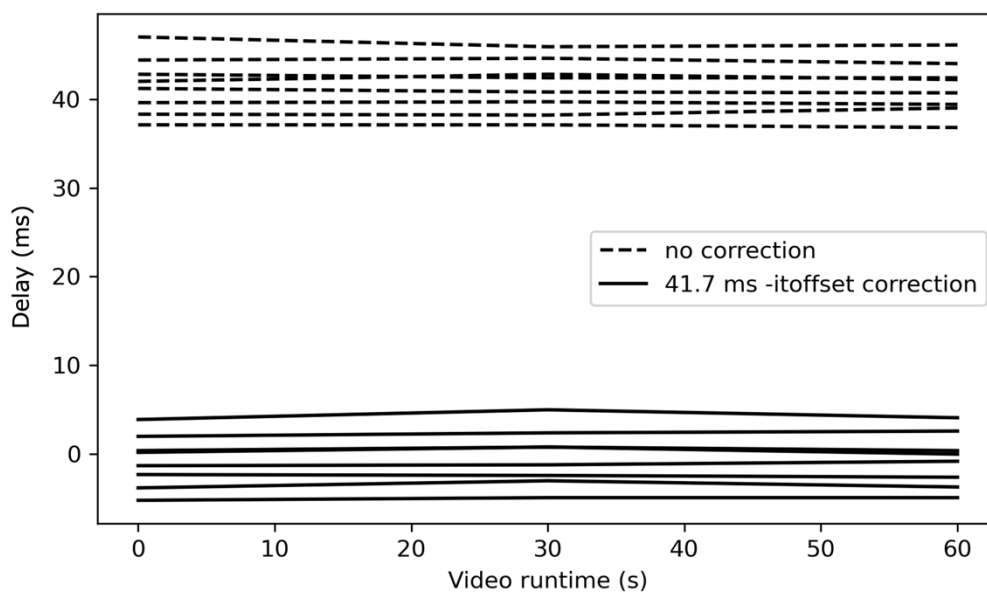


Figure 4 — Delay of the audio signal with respect the visual stimulus over the duration of a synchronisation video, for several trials. Before correcting for the delay in the presentation software (dashed), the average delay is 41.7 ms. After correction (solid) we obtain a negligible delay of 0.6 ms.

¹ Available at <https://pypi.org/project/python-vlc/>

² Available at <https://ffmpeg.org>

Comprehension Tests

Participants

To validate our multiple-choice comprehension questions, we conducted a behavioural experiment with seventeen native English speakers, ten of them female, with self-reported normal or corrected-to-normal vision and healthy hearing. The participants were between 18 and 29 years of age, with a mean age of 23 years. All participants were right-handed and had no history of mental health problems, severe head injury or neurological disorders. Before starting the experiment, participants gave informed consent. The experimental protocol was approved by the Imperial College Research Ethics Committee.

Stimuli Presentation

The experiment took place in an acoustically and electrically insulated room (IAC Acoustics, United Kingdom). The same equipment and settings employed in the audiovisual synchronisation experiment were used to control the audiovisual presentation and data acquisition.

We considered three types of stimuli. The first was audio-only speech in stationary, speech-shaped background noise. The second condition was audiovisual speech, also in background noise. In both conditions the speech was presented in a constant level of background noise, at a signal to noise ratio of -2 dB. The third condition was visual-only speech, with no sound presented. Subjects listened to 54 passages from the corpus and the three conditions were randomised between the different passages. Between each passage, the participants were tasked with answering the AVbooks comprehension questions and then to read the summary before proceeding to the next task. A speech comprehension score was computed as the average percentage of correct answers for each of the three conditions.

Results

The comprehension scores (Figure 5) in the audio-only condition were found to be $55.3\% \pm 3.5\%$ (mean and standard error of the mean) while they were $60.6\% \pm 2.8\%$ in the audiovisual condition. The difference between the scores in the two conditions was not significant ($p = 0.18$, $t = 1.40$, paired Student's t-test with Benjamini-Hochberg FDR correction). On the other hand, subjects scored significantly lower in the video only condition ($34.1\% \pm 2.9\%$) than when presented with the audio signal only or with the audiovisual stimuli ($p = 1.1 \times 10^{-5}$, $t = 6.51$ and $p = 3.5 \times 10^{-6}$, $t = 7.56$ respectively, FDR corrected paired Student's t-tests). Participants did not score better than the chance level (33%) in the video only condition ($p = 0.72$, $t = 0.37$, one sample Student's t-test).

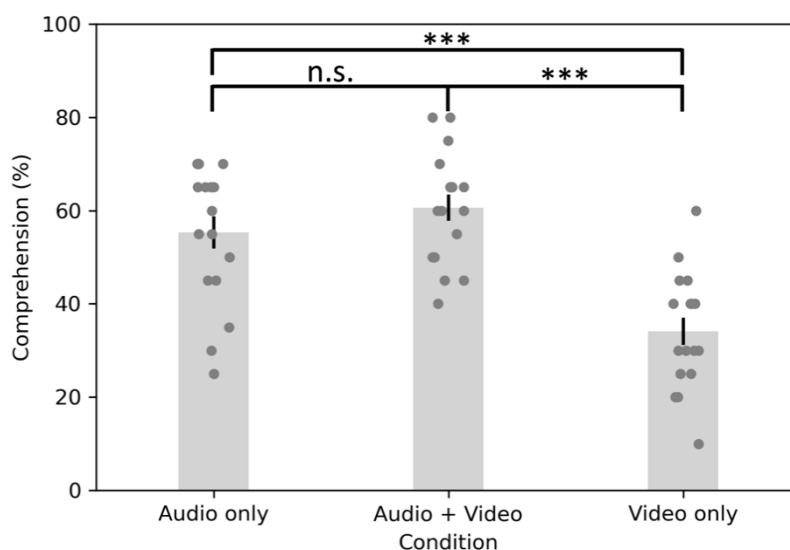


Figure 5 — Speech comprehension scores for the AVbook in background noise without a visual signal (audio only), for the audiovisual stimulus, as well as for the video only. Error bars represent the standard error of the mean, and the dark grey points show the average score per participant.

Conclusion

A narrated, continuous speech audiovisual corpus, AVbook, was collected to support the investigation of audiovisual speech perception through brain imaging and behavioural studies. Although the corpus is limited to two speakers, the high frame rate may make it an attractive option for developing and benchmarking computer algorithms of audiovisual speech processing as well.

The corpus passages were chosen to yield high attention and engagement of experimental subjects. In case the passages were not fully comprehensible to participants, for instance due to the inclusion of background noise, a short written summary was produced for each video segment. It is hoped that providing this summary text before presenting the next passage will help to keep a subject's engagement at a high level.

We developed comprehension questions that accompany the AVbook corpus, and tested these in a behavioural experiment. Our finding that participants scored well above chance when they could hear the speech material, but that their performance was at chance level when they could only lip-read suggests that the question set is fit for the purpose of checking for attention. On the other hand, our behavioural experiment revealed no significant difference between the audio-only and the audiovisual conditions. This presumably reflected that the questions were designed to test participants' understanding of the stories and attention to it, and not for a highly accurate quantification of comprehension. The latter would benefit from single semantically unpredictable sentences that are assessed for the comprehension of key words. Furthermore, a difference may emerge in or across other listening conditions (SNR, noise type) and with future work aiming to characterise the Q&A set while disambiguating working memory effects and attention effects from comprehension, through careful experimental design.

We also developed a solution, an electronic clapperboard, to temporally align the audio and video stimuli precisely and reliably. Our synchronisation experiment, aided by the videos of the electronic clapperboard, yielded a negligible mean latency that remained consistent within trials and jittered by less than one frame (8.1 ms vs 8.3 ms) across trials.

Due to the precise audiovisual synchronization, the AVbook corpus can allow the neuroscience community to extend work done in the audio-only modality with traditional audio-book corpora to include the visual domain. Such studies could tackle the effect of audiovisual speech on neural responses related to attention (O’Sullivan et al., 2015) or linguistic factors such as surprisal (Weissbart et al., 2020). Furthermore, a consistent use of audiovisual speech material across geographies, equipment and experimental paradigms may help compare and reproduce results.

Further, due to its consistent framing and high number of frames, the AVbook corpus may be an attractive option for computational studies and therefore support the use of common material in language, speech perception and automatic speech recognition work, a result largely achieved by the GRID and VIDTIMIT corpora in the context of short sentences.

The complete corpus, phone-level segmentation files, synchronization material, teleprompter scripts, comprehension questions, summary texts, and raw video recordings are available to download for research use at <https://zenodo.org/record/7387046>.

Acknowledgements

This research was supported by the Royal British Legion Centre for Blast Injury Studies, EPSRC grants EP/M026728/1 and EP/R032602/1, as well as by the U.S. Army through project 71931-LS-INT.

References

- Benezeth, Y., & Bachman, G. (2011). BL-Database: A French audiovisual database for speech driven lip animation systems. <http://hal.inria.fr/inria-00614761/>.
- Bernstein, L. E., Auer, E. T. Jr., and Takayanagi, S. (2004). Auditory speech detection in noise enhanced by lipreading. *Speech and Communication*, 44, 5–18. doi: 10.1016/j.specom.2004.10.011
- Brookes, M. (2022). “Speech processing toolbox for MATLAB,” available from: <https://github.com/ImperialCollegeLondon/sap-voicebox>, commit: 38061fe.
- Brown, V.A., Hedayati, M., Zanger, A., Mayn, S., Ray, L., et al. (2018) What accounts for individual differences in susceptibility to the McGurk effect?. *PLOS ONE* 13(11): e0207160. doi: 10.1371/journal.pone.0207160
- Chitu, A.G., and Rothkrantz, L.J.M. (2007). Building a data corpus for audio-visual speech recognition. In L. J. M. Rothkrantz, & C. A. P. G. van der Mast (Eds.), *Euromedia 2007* (pp. 88-92). Eurosis-ETI.
- Cooke, M., Barker, J., Cunningham, S., & Shao, X. (2006). An audio-visual corpus for speech perception and automatic speech recognition. *The Journal of the Acoustical Society of America*, 120 (5 Pt 1), 2421–2424. doi: 10.1121/1.2229005.
- Crosse, M. J., & Lalor, E. C. (2014). Retraction. *Journal of neurophysiology*, 112(10), 2667. doi: 10.1152/jn.z9k-2710-retr.2014
- Crosse, M.J., Di Liberto, G.M., Bednar, A., Lalor, E.C. (2016). The Multivariate Temporal Response Function (mTRF) Toolbox: A MATLAB Toolbox for Relating Neural Signals to Continuous Stimuli. *Frontiers in Human Neuroscience*. 2016;10:604-604. doi: 10.3389/fnhum.2016.00604
- Grant, K. W., and Seitz, P. F. (2000). The use of visible speech cues for improving auditory detection of spoken sentences. *The Journal of the Acoustical Society of America*, 108, 1197–1208. doi: 10.1121/1.1288668
- ITU-T (2011). “ITU-T P.56: Objective measurement of active speech level (12/2011),” International Telecommunication Union Recommendation E 37305, Series P: Terminals and Subjective and Objective Assessment Methods – Objective Measuring Apparatus, Geneva, Switzerland.
- Matthews, I., Cootes, T.F., Bangham, J.A., Cox, S., Harvey, R. (2002). "Extraction of visual features for lipreading," in *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 2, pp. 198-213, doi: 10.1109/34.982900.

- Lee, B., Hasegawa-johnson, M., Goudeseune, C., Kamdar, S., Borys, S., Liu, M., & Huang, T. (2004). AVICAR : Audio-Visual Speech Corpus in a Car Environment. *8th International Conference on Spoken Language Processing* pp. 8–11.
- Luck, S. J. (2014). *An Introduction to the Event-Related Potential Technique*. Cambridge, MA: MIT press.
- Marmarelis, V. Z. (2004). *Nonlinear Dynamic Modeling of Physiological Systems*. Hoboken, NJ: John Wiley & Sons.
- McGurk, H., MacDonald, J. (1976). Hearing lips and seeing voices. *Nature* 264, 746–748. doi: 10.1038/264746a0
- Messer, K., Matas, J., Kittler, J., & Jonsson, K. (1999). XM2VTSDB: The extended M2VTS database. In *Second International Conference on Audio and Video-based Biometric Person Authentication* (pp. 72–77).
- Movellan, J. R. (1995). "TULIPS1 database" in *Visual Speech Recognition with Stochastic Networks in Advances in Neural Information Processing Systems*, Cambridge:MIT Press, vol. 7.
- Ding, N., Simon, J.Z., (2014). Cortical entrainment to continuous speech: Functional roles and interpretations. *Frontiers in Human Neuroscience*. 8, 311
- Ding, N., Simon, J.Z. (2012). Emergence of neural encoding of auditory objects while listening to competing speakers. *Proceedings of the National Academy of Sciences of the United States of America* 109, 11854–11859.
- O'Sullivan, J.A., Power, A.J., Mesgarani, N., Rajaram, S., Foxe, J.J., Shinn-Cunningham, B.G., Slaney, M., Shamma, S.A., Lalor, E.C. (2015). Attentional Selection in a Cocktail Party Environment Can Be Decoded from Single-Trial EEG, *Cerebral Cortex*, Volume 25, Issue 7, July 2015, Pages 1697–1706, doi: 10.1093/cercor/bht355
- O'Sullivan, A., Crosse, M.J., Di Liberto, G.M., Lalor, E.C. (2017). Visual Cortical Entrainment to Motion and Categorical Speech Features during Silent Lipreading. *Frontiers in Human Neuroscience*. doi: 10.3389/fnhum.2016.00679
- Reisberg, D., Mclean, J., and Goldfield, A. (1987). "Easy to hear but hard to understand: a lip-reading advantage with intact auditory stimuli," in *The Psychology of Lip-Reading*, eds B. Dodd and R. Campbell (Hillsdale: Lawrence Erlbaum Associates), 97–114.
- 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. *Cerebral Cortex* 17, 1147–1153. doi: 10.1093/cercor/bhl024
- Sanderson, C., Lovell, B.C. (2009). Multi-Region Probabilistic Histograms for Robust and Scalable Identity Inference. In: *Tistarelli, M., Nixon, M.S. (eds) Advances in Biometrics. ICB 2009. Lecture Notes in Computer Science*, vol 5558. Springer, Berlin, Heidelberg. doi: 10.1007/978-3-642-01793-3_21.
- Schultz B.G., Biau E., Kotz S.A. (2020). An open-source toolbox for measuring dynamic video framerates and synchronizing video stimuli with neural and behavioral responses. *Journal of Neuroscience Methods*, 343 (2020), Article 108830
- Sonkusare, S., Breakspear, M., & Guo, C. (2019). Naturalistic Stimuli in Neuroscience: Critically Acclaimed. *Trends in cognitive sciences*, 23(8), 699–714. <https://doi.org/10.1016/j.tics.2019.05.004>
- Sumby, W. H., and Pollack, I. (1954). Visual contribution to speech intelligibility in noise. *The Journal of the Acoustical Society of America*, 26, 212–215. doi: 10.1121/1.1907309
- Supasorn, S., Seitz, S.M., and Kemelmacher-Shlizerman, I. (2017). Synthesizing Obama: learning lip sync from audio. *Association for Computing Machinery Transactions on Graphics*. 36, 4, Article 95, 13 pages. doi: 10.1145/3072959.3073640
- Weissbart, H., Kandylaki, K.D., Reichenbach; T. (2020). Cortical Tracking of Surprisal during Continuous Speech Comprehension. *Journal of Cognitive Neuroscience* 2020; 32 (1): 155–166. doi: 10.1162/jocn_a_01467

Comments

Comments from Baptiste Bouvier:

I appreciated reading this clear and promising study. In particular, I understood that what is innovative and particularly interesting in the dataset the authors propose is:

The presentation of complex and large audiovisual stimuli (whole texts read continuously) ;

The fine precision in the quality of both audio and visual stimuli, and the synchronization between these audio and visual stimuli enabled by recordings at a high frame rate.

Congrats for all the work of acquisition and processing of stimuli, which is rich and opens wide perspectives for future neuroscience studies.

Author's Response:

Many thanks for your review and questions.

Here are my comments:

Page 3, line 90: "The latter action was taken to aid the speakers in producing a 90 neutral tone without overly salient parts." What do the authors mean by 'salient parts' ? And why removing them from the text ?

Author's Response:

Our mention of salient parts refers to direct speech and passages intended to make the reader laugh or giggle through play on words. These were removed to reduce the risk of EEG artifacts.

About the comprehension test

The purpose of this test is to validate the multiple choice questionnaire. The authors proposed three conditions, one with the audio-visual stimulus in noise, one with the audio stimulus alone in noise, and one with the visual stimulus without noise. I have a few comments about the choice of these conditions.

In particular, what is the purpose of the condition with the visual stimulus alone ? I understand that the authors wanted to check that the rate of correct answers to the questionnaires without having listened to the speech is the same as chance. But why not just have the questionnaire done without presenting the stimuli at all ? It seems that presenting the visual stimulus alone only tests the participants' ability to lip-read there.

How can the results of this comprehension test be interpreted in the perspective of the questionnaire evaluation ? As it stands, it seems to me that the conclusions of this test are :

The quiz is not feasible at a better-than-lucky rate if one has not heard the speeches.

The correct response rate when one has heard the speech is just under or around 60%, which is better than chance but could be higher : isn't the questionnaire too complicated ? I understand however that it has to be under 100% so that the participants have to concentrate to the maximum.

Adding the visual stimulus on top of the audio stimulus does not improve the scores : what does this result add to the evaluation of the quality of the questionnaire ?

Thus, I wonder how these results really validate the questionnaire.

Author's Response:

We now clarify that the purpose of the Q&A set is generally to “check[ing] for attention” (line 255). The authors agree that it is important to further characterise (and optimise) the Q&A set, perhaps in different listening conditions and across different populations, and with experimental designs capable of disambiguating comprehension from attention and working memory. We also now discuss (lines 260-263):

Furthermore, a difference may emerge in or across other listening conditions (SNR, noise type) and with future work aiming to characterise the Q&A set while disambiguating working memory effects and attention effects from comprehension, through careful experimental design.

Moreover, the authors wrote on page 7, line 242, that "Providing this summary text before presenting the next passage will help to keep a subject's engagement at a high level", but this doesn't seem obvious from the study. To validate the summaries, it would be needed to compare a condition where participants have access to the summaries of each passage after answering the questions in that passage and one where they do not.

Author's Response:

We now caution that “It is hoped that providing...” (line 250).

In summary, for this test of questionnaire validation (which could be combined with passage summary validation), why not instead compare questionnaire scores under the following conditions:

Audio-visual stimulus with summaries between passages

Audio-visual stimulus without summaries

No stimulus

Author's Response:

Many thanks for this proposal. Further work in these directions will complement the corpus as provided and is a welcome contribution.

Comments from Pedro Lladó:

This study presents a high frame rate corpus of audiovisual speech, accurately time-aligned for multimodal speech perception. The corpus consists of passages of a book read by two speakers (1h 50' each). A comprehension test based on multiple choice questions was proposed and tested in three conditions: audio only, video only and audio+video. The results showed that comprehension in the video only condition was about chance level. The results obtained in the audio only and audio+video conditions were reasonably similar, both well above chance level.

In my opinion, the study presents interesting contributions to the field:

The corpus has been accurately collected, with controlled valuable features for multimodal perception research.

A method to test comprehension of longer speech fragments has been proposed as opposed to existent short sentences speech comprehension tests.

The method to control the synchronization seems good and easy to implement, which should be included in future research.

Despite I don't have any major concern about the work, I would like to make some questions:

Author's Response:

Many thanks for your review and questions.

Line 189: Could you explain how was the delay corrected? I assumed that the average delay would be zero after correction, but it is not.

Author's Response:

The delay was corrected using ffmpeg and its "itoffset" flag. This is already stated in line 183. The residual error is due to jitter and uncertainty in measurement.

Related to the previous question: could you provide some info/references about what is an acceptable delay for multimodal speech perception? Same for the jitter. It would help the reader to understand if the correction is good enough.

Author's Response:

Temporal window of integration in auditory-visual speech perception (2007) V. van Wassenhove, K. W. Grant, D. Poeppel

About the multiple-choice comprehension questions: Did you base it on any test that was done before? Why this test was chosen? I understand that this is proposed as opposed to shorter sentences comprehension tests, but you could make it explicit in the manuscript to convince the reader that the test was necessary and appropriate.

Author's Response:

We now clarify that the purpose of the Q&A set is generally to "check[ing] for attention" (line 255). We also now discuss (lines 260-263):

Furthermore, a difference may emerge in or across other listening conditions (SNR, noise type) and with future work aiming to characterise the Q&A set while disambiguating working memory effects and attention effects from comprehension, through careful experimental design.

Further work in these directions will complement the corpus as provided and is a welcome contribution.

The results of the audio only and audio + video conditions are on average 55.3% and 60.6% respectively. What is the reason for it being so "far" from perfect comprehension? Is it because of intelligibility? Would you expect perfect scores without the noise? In your opinion, would the SNR have an effect on the differences between these two conditions? I think there is room for interesting discussion when analyzing the results.

Author's Response:

Due to the (presumably complex) interaction of working memory, comprehension and attention constraints, we would not expect perfect responses to the questions.

I couldn't find the link in the manuscript and I couldn't find the data on zenodo. Could you provide the link in the manuscript?

Author's Response:

We now state: The complete corpus, phone-level segmentation files [NEW], synchronization material, teleprompter scripts, comprehension questions, summary texts, and raw video recordings are available to download for research use at <https://zenodo.org/record/7387046>.

Comments from Mengchao Zhang:

The current paper develops an audio-visual speech corpus that implemented continuous speech materials from the book 'Endurance: Shackleton's Incredible Voyage to the Antarctic'. Continuous speech materials are advantageous due to its closeness to speech in natural communication. The study also improved the recording of the material so that the auditory and visual signals can be presented with high visual resolution and precise AV synchrony.

The entire recording and rendering process is beneficial for the broader research community to create high-precision AV stimuli. The paper is well written and very informative. I have some questions in terms of the choice of the source material and the interpretation of the behavioral task.

I am curious about the rationale of choosing this book. As participants may receive the book to different extents, different levels of enjoyment or familiarity with the book may introduce unwanted variabilities to tasks like comprehension or speech recognition. Is there a reason that a more neutral source material was not selected?

Author's Response:

Many thanks for your review and comments. The rationale behind the choice of book was indeed to reduce unwanted variabilities due to familiarity and political or cultural tension to the largest extent possible, while maximising enjoyment of the storytelling to improve participant attention in long M/EEG studies. Although subjective, in the authors' opinion the chosen text is generally enjoyable and politically neutral, unlike a commonly-employed corpus of political speeches and addresses by a well-known public figure.

Line 244-252: One could interpret that comprehension under audio-only and audiovisual conditions are similar because participants used the auditory signals to comprehend the story. It is not entirely clear how the comprehension task assessed attention to the story, not comprehension. More details on the rationales and the procedure of comprehension task would be helpful.

Author's Response:

The authors agree that it is important to further characterise (and optimise) the Q&A set, perhaps in different listening conditions and across different populations, and with experimental designs capable of disambiguating comprehension from attention and working memory. Further work in these directions will complement the corpus as provided and is a welcome contribution. We now state (lines 260-263):

Furthermore, a difference may emerge in or across other listening conditions (SNR, noise type) and with future work aiming to characterise the Q&A set while disambiguating working memory effects and attention effects from comprehension, through careful experimental design.

Comments from Deniz Başkent:

This study presents an AV book recording where a lot of work has been done to control for A-V asynchrony. As someone who had to deal with this in previous research, I appreciated this effort very much and the materials will be a valuable tool to researchers.

When we wanted to make AV recordings it turned out to be technically difficult, not only for the AV synchrony issues, but also the rooms where we could make excellent quality A recordings were not optimal for V recordings and vice versa (at least at our institution). Authors explain the V quality in the paper but A part seems a bit less explained (only sampling and bit rate). I understood there was no soundcard in the A recording setup, and I was also not sure how the acoustics environment of the recording studio was, and from the demos during presentation it was also not very easy to evaluate this. Since the A quality could have an effect of speech comprehension in noise, it would be useful to add some details on this also. On the lines of, how the authors ensured this quality and were the A recordings free from any potential noises from room acoustics or camera auxiliary connection.

Author's Response:

Many thanks for your appraisal and comments. To clarify, although no separate audio card was employed in the recording, care was taken to use high-quality devices capable and a clip-on microphone to improve SNR. The SNR is now reported in the text and we now state (lines 123-131):

No clipping was observed in the audio and no further processing was performed.

The audio was also separately exported as a WAVE file. An a priori signal-to-noise ratio (SNR) was estimated using a voice activity detection (VAD) script (Brookes, 2022) following the ITU-T P.56 standard (ITU-T, 2011). The average SNR was found to be 34.0 ± 2.9 dB and 44.5 ± 7.1 dB for the female and male speakers respectively (mean and standard deviation). Similar results were obtained when employing a custom-tuned VAD script or by calculating the SNR by comparing silent video segments which were cut out of the final corpus (less than 1 dB difference in both cases for both speakers).

Up-to-date comments can be found on [PubPeer](#).