# Quantifying the individual auditory and visual brain response in 7- month-old infants watching a brief cartoon movie

Sarah Jessen[1][#], Lorenz Fiedler[2], Thomas F. Münte[1], Jonas Obleser[2]


[1]: Department of Neurology, University of Lübeck, Lübeck, Germany

[2]: Department of Psychology, University of Lübeck, Lübeck, Germany




[#] Lead contact:

Dr Sarah Jessen, Department of Neurology, University of Lübeck, Ratzeburger Allee 160, 23562 Lübeck, Germany, Email: sarah.jessen@neuro.uni-luebeck.de

Phone: +49 451 3101 7449

Running title: Quantifying the individual infant brain response to real-life stimuli

Number of pages: 24

Number of figures: 8

Number of words (manuscript): 6240

Number of words (abstract): 245

**Keywords**: EEG; audiovisual; forward encoding models; temporal response function; ecologically valid stimuli; developmental neuroscience

26    ABSTRACT

27    Electroencephalography (EEG) continues to be the most popular method to investigate
28    cognitive brain mechanisms in young children and infants. Most infant studies rely on the
29    well-established and easy-to-use event-related brain potential (ERP). As a severe
30    disadvantage, ERP computation requires a large number of repetitions of items from the
31    same stimulus-category, compromising both ERPs' reliability and their ecological validity
32    in infant research. We here explore a way to investigate infant continuous EEG responses
33    to an ongoing, engaging signal (i.e., "neural tracking") by using multivariate temporal
34    response functions (mTRFs), an approach increasingly popular in adult-EEG research.
35    N=52 infants watched a 5-min episode of an age-appropriate cartoon while the EEG signal
36    was recorded. We estimated and validated forward encoding models of auditory-envelope
37    and visual-motion features. We compared individual and group-based ('generic') models of
38    the infant brain response to comparison data from N=28 adults. The generic model yielded
39    clearly defined response functions for both, the auditory and the motion regressor.
40    Importantly, this response profile was present also on an individual level, albeit with lower
41    precision of the estimate but above-chance predictive accuracy for the modelled individual
42    brain responses. In sum, we demonstrate that mTRFs are a feasible way of analyzing
43    continuous EEG responses in infants. We observe robust response estimates both across
44    and within participants from only five minutes of recorded EEG signal. Our results open
45    ways for incorporating more engaging and more ecologically valid stimulus materials when
46    probing cognitive, perceptual, and affective processes in infants and young children.

47

2

48 INTRODUCTION

49 Neuroimaging studies in healthy human infants are subject to severe constraints, as

50 participants cannot follow verbal instructions, show generally short attention spans, and

51 overall tend to be not very cooperative. As functional magnetic resonance imaging (fMRI)

52 studies are difficult to realize in infants (Ellis & Turk-Browne, 2018),

53 electroencephalography (EEG) continues to be the most popular method to investigate

54 cognitive brain mechanisms in very young children and infants.

55 To analyze the EEG signal, most studies in infants rely on the use of event-related brain

56 potentials (ERPs). Accordingly, most infant EEG paradigms have been optimized for the

57 computation of ERPs: This method necessitates that a few, carefully selected stimulus

58 conditions are repeated multiple times to elicit and average a stereotypical brain response

59 (i.e., an ERP) that can then be compared between conditions or between individuals. This

60 leads to experimental designs that are often (a) highly unnatural and (b) have difficulties

61 capturing the infants' attention for more than a few minutes.

62 However, in recent years and with the advent of modern computational possibilities,

63 several new approaches to analyze EEG data have become available in adult EEG research.

64 One such approach is the so-called "neural tracking", which seeks to compute and assess

65 the relationship between the recorded EEG signal and an ongoing stimulus signal. The key

66 ideas here are, first, naturally varying, non-repetitive stimuli, often movies (Bartels, Zeki,

67 & Logothetis, 2008; Hasson, Nir, Levy, Fuhrmann, & Malach, 2004; Nishimoto et al.,

68 2011) or naturally spoken conversation (Broderick, Anderson, Di Liberto, Crosse, & Lalor,

69 2018; Ding & Simon, 2013; Fiedler, Wöstmann, Herbst, & Obleser, 2019), which have

70 higher ecological validity and arguably engage the participant qualitatively differently than

71 artificial, isolated stimuli (Hamilton & Huth, 2018; Huk, Bonnen, & He, 2018; Matusz,

72 Dikker, Huth, & Perrodin, 2018). Second, a mathematical framework (usually a variant of

73 the general linear model) that allows to either "reconstruct" features of such a natural

74 stimulus based on the ongoing brain response (so-called backward or decoding models), or

75 to "predict" the measured ongoing brain response from features of the stimulus (so-called

76 forward or encoding models; Dayan & Abbott, 2001; Naselaris, Kay, Nishimoto, &

77 Gallant, 2011).

78 While the use of these advanced EEG analysis approaches has become rapidly mainstream

79 in non-human and adult human neuroscientific research, it is still rare in infant research.

80 This is unfortunate, since they not only have yielded important new insights in adult

81 research and are likely to offer the same potential in infant studies, but they may even

3

provide higher gains in infancy research, which suffers from notoriously low data quality and quantity. It may for instance reduce attrition rates, as experimental designs can be optimized to be highly engaging for infant participants. Rather than presenting hundreds of repetitions of very similar stimuli, which raises the additional challenge of keeping a non-cooperative participant attending to the screen, participants can be presented with constantly changing, engaging videos in which stimuli are embedded.

Importantly, as in adult work, infant brain research has seen an increased interest in the use of naturalistic settings over the past years. Recent research has for instance demonstrated the feasibility of investigating interpersonal neural coupling in adult-infant-interactions (Leong et al., 2017) or the use of oscillatory brain responses in analyzing responses to dynamic social information (Jones, Venema, Lowy, Earl, & Webb, 2015). While dynamic, naturalistic settings and experimental paradigms yield important new insights into how brains behave and interact in real life rather than an abstract laboratory setting, they inherently pose the additional challenge of hard-to-predict and highly variant sensory input. Being able to directly relate a constantly changing input to ongoing brain responses would therefore also be crucial for the analysis of state-of-the-art ecologically valid experimental designs.

One particularly promising approach to do so is the use of multivariate temporal response functions (mTRFs), which offer a mathematically simple way to link ongoing, continuous environmental signals to simultaneously recorded brain responses. In adults, mTRFs have successfully been used to track the processing of ongoing speech (e.g., Fiedler et al., 2019) as well as ongoing and naturalistic visual input (O'Sullivan, Crosse, Di Liberto, & Lalor, 2017). Furthermore, Kalashnikova et al. (2018) used mTRFs in infants to analyze the processing of ongoing auditory speech signals, reporting a stronger cortical tracking for infant-directed compared to adult-directed speech (Kalashnikova, Peter, Di Liberto, Lalor, & Burnham, 2018).

We here demonstrate the feasibility and utility of a forward encoding modelling combined with non-repetitive complex multisensory stimulation in an infant population. We presented 7-month-old infants with a 4'48'' long age-appropriate cartoon (one episode of the cartoon-show *Peppa Pig*) while recording the EEG. We focused our analysis on the processing of three low-level physical stimulus parameters; the auditory envelope, the motion content, and luminance. All three parameters have been amply investigated in both infants and adults and are known to elicit reliable ERP responses.

4

115  The auditory ERP response typically consists of a frontocentral P1–N1–P2–N2 sequence

116  of responses, which can be clearly observed in adults and emerges in infancy and early

117  childhood (see e.g., Wunderlich & Cone-Wesson, 2006, for a review). Compared to adults,

118  infants tend to show a much less pronounced P1–N1 response, and the overall response is

119  dominated by a broad P2 response (Wunderlich, Cone-Wesson, & Shepherd, 2006).

120  The infant visual ERP to complex stimuli such as objects and faces comprises three main

121  components; the Pb, the Nc, and the Slow Wave (Webb, Long, & Nelson, 2005). In

122  particular, the Nc response, a frontocentral negativity typically observed between 400 and

123  800 ms after stimulus onset often linked to the allocation of attention has been amply

124  investigated (de Haan, Johnson, & Halit, 2003; Reynolds & Guy, 2012).

125  If we were successful in estimating auditory and visual brain responses using a forward

126  encoding model approach, we expect response functions comparable to classical evoked

127  brain responses. Furthermore, since the combined use of auditory and visual regressors

128  provides more information compared to the use of either regressor alone, we expected a

129  more consistent and reliable response function when using auditory and visual regressors

130  in one model.

131  Finally, while it is common in adult studies using mTRFs to compute individual response

132  functions based on a subset of the available data, due to the limited amount of data available

133  in the infant cohort we aimed to explore the potential benefit from relying on a "generic"

134  response function (Di Liberto & Lalor, 2017), that is, an average response function

135  computed across participants. Hence, we computed an averaged response function over $n$–

136  1 participants and used this response function to model responses in the $n$th participant (i.e.

137  leave-one-out cross validation). We directly contrasted results obtained with these two

138  approaches on the present data set.

139

140  METHODS

141  *Infant participants*. Fifty-two 7-month-old infants were included in the final sample (age:

142  213 ± 8 days [mean ± standard deviation (SD)]; range: 200-225, 24 female). Not untypical

143  for infant studies (Stets, Stahl, & Reid, 2012), an additional 39 infants had been tested but

144  could not be included in the final sample. Note also that directly prior to the experiment

145  reported here, infants had already participated in a 5–10-minute-long ERP experiment on

146  visual emotion perception (see below), further contributing to the drop-out rate since

147  infants often became fussy or tired after the first experiment. In detail, infants were

5

148   excluded because they did not watch the complete video (n=24); were too fussy to watch
149   the video at all (n=10); did not contribute at least 100 s of artifact-free data (n=3); had
150   potential neurological problems (n=1); or because of technical problems during the
151   recording (n=1).

152   All infants were recruited via the maternity ward at the local hospital (Universitätsklinikum
153   Schleswig-Holstein); were born full-term (38–42 weeks gestational age); had a birth weight
154   of at least 2500 g; and had no known neurological deficits. The study was conducted
155   according to the declaration of Helsinki, approved by the ethics committee at the University
156   of Luebeck, and parents provided written informed consent.

157   *Adult reference sample.* In addition, we collected data from a reference sample of n = 33
158   adult participants. Data from n=5 were excluded due to technical difficulties during the
159   recording (n=2) or failure to contribute at least 100 s of artifact-free data (n=3), leading to
160   a final sample of n=28 (mean age: 50 years; range: 21–69, 16 female).

161   *Stimulus.* As stimulus material we used one episode (duration 4'48'', that is, 269 s or 6451
162   frames) of the cartoon show Peppa Pig ("Peppa Pig–The new car"), an age-appropriate
163   cartoon featuring a family of pigs and their daily life. Sound and visual parameters were
164   not manipulated in any way.

165   *Procedure–Infants.* After arrival in the laboratory, parents and infant were familiarized
166   with the environment and parents were informed about the study and signed a consent form.
167   The EEG recording was prepared while the infant was sitting on his/her parent's lap. For
168   recording, we used an elastic cap (BrainCap, Easycap GmbH) in which 27 Ag/AgCl-
169   electrodes were mounted according to the international 10-20-system. An additional
170   electrode was attached below the infant's right eye to record the electrooculogram. The
171   EEG signal was recorded with a sampling rate of 250 Hz using a BrainAmp amplifier and
172   the BrainVision Recorder software (both Brain Products).

173   For the EEG recording, the infant was sitting in an age appropriate car seat (Maxi Cosi
174   Pebble) positioned on the floor. As part of a larger study, a t-shirt was positioned over the
175   chest area of the infants. The t-shirt had either previously been worn by the infant's mother
176   (n=19) or by the mother of a different same-aged infant (n=14) or had not been worn before
177   (n=19). This modulation was not of main interest to the present study and will not be
178   analyzed or reported here in further detail.

179   In front of the infant (approximately 60 cm from the infant's feet), a 24-inch monitor with
180   a refresh rate of 60 Hz was positioned at a height of about 40 cm (bottom edge of the

6

screen). Left and right of the monitor, loudspeakers (Logitech X-140) were positioned and set to a comfortable level of loudness. When the infant was attending to the screen, the video was started and played without interruption until the end of the episode. The parent was seated approximately 1.5 m behind the infant and was instructed not to interact with the infant during the video. In case the infant became too fussy and started crying during the video, the video was aborted and the infant was excluded from further analysis.

Before this video presentation, infants had been presented with a series of photographs displaying happy and fearful facial expressions as part of the larger, maternal-odor study. Again, the results of this part of the study will not be further analyzed here.

*Procedure–Adults.* Adult participants were presented with the same "Peppa Pig" movie after they had already participated in one of several unrelated EEG studies. They were informed about the study and signed a consent form. For recording the EEG signal, we used 64 Ag/AgCl active scalp electrodes positioned in an elastic cap according to the international 10-20-system. The EEG signal was recorded with a sampling rate of 1000 Hz using an ActiChamp amplifier and the BrainVision Recorder software (Brain Products).

Adult participants sat in a soundproof and electrically shielded chamber (Desone) in a comfortable chair approximately 1 m away from a 24-inch monitor with a refresh rate of 60 Hz on which the video was presented. Sound was presented from the same loudspeaker models used in the infant study, also positioned left and right to the screen (Logitech X-140).

*Analysis.* Unless noted otherwise, the analysis protocol was identical for infant and adult data. We analyzed the data using Matlab 2013b (The MathWorks, Inc., Natick, MA), the Matlab toolbox FieldTrip (Oostenveld, Fries, Maris, & Schoffelen, 2011), and the multivariate temporal response function (MTRF) toolbox (Crosse, Di Liberto, Bednar, & Lalor, 2016).

*Preprocessing.* The data were referenced to the average of all electrodes (mean reference), filtered using a 100-Hz-lowpass and a 1-Hz-highpass filter, and segmented into 1-sec-epochs. To detect epochs contaminated by artifacts, the standard deviation was computed in a sliding window of 200 msec. If the standard deviation exceeds 80 mV at any electrode, the entire epoch was discarded, and if less than 100 artifact-free epochs remained, the participant was excluded from further analysis. An independent component analysis (ICA) was computed on the remaining epochs. Components were inspected visually by a trained coder (S.J.) and rejected if classified as artefactual (infants: $5 \pm 2$ components per participant [mean $\pm$ SD], range 1–10; adults: $26 \pm 5$, range 11–36). A 1–10 Hz bandpass

215    filter was applied to the cleaned data. Adult data were downsampled to the infant-data

216    sampling frequency of 250 Hz.

217

218        *Extraction of stimulus regressors.* Regressors characterizing motion, luminance,

219    and the sound envelope were extracted from the stimulus video. Exemplary excerpts of

220    audio, luminance, and motion regressors are shown in Fig 1B.

221        To compute a regressor of average luminance across all pixels, the weighted sum

222    of the rgb values for each frame was computed using Matlab (Bartels et al., 2008).

223        To compute a regressor of average motion across all pixels, each video frame was

224    converted to grey-scale, and the difference between two consecutive frames was computed.

225    Then, the mean across all pixels for which this difference was larger than 10 (to account

226    for random noise, see e.g. Jessen & Kotz, 2011; Pichon, de Gelder, & Grèzes, 2009) was

227    computed.

228        To compute a regressor of sound envelope, the audio soundtrack of the video was

229    extracted and submitted to the NSL toolbox, an established preprocessing pipeline

230    emulating important stages of auditory peripheral and subcortical processing (Ru, 2001).

231    The output of this toolbox resulted in a representation containing band-specific envelopes

232    of 128 frequency bands of uniform width on the logarithmic scale with center frequencies

233    logarithmically spaced between 0.1 and 4 kHz. To obtain the broadband temporal envelope

234    of the audio soundtrack, these band-specific envelopes were then summed up across all

235    frequencies to obtain one temporal envelope. Following earlier own and others'

236    approaches, we used the first derivative of the half-wave rectified envelope as the final

237    audio regressor (for details see Fiedler et al., 2017). The result is a pulse-train-like series

238    of peaks where, across frequency bands, the acoustic energy rises most steeply, reflecting

239    "acoustic edges" such as syllable onsets.
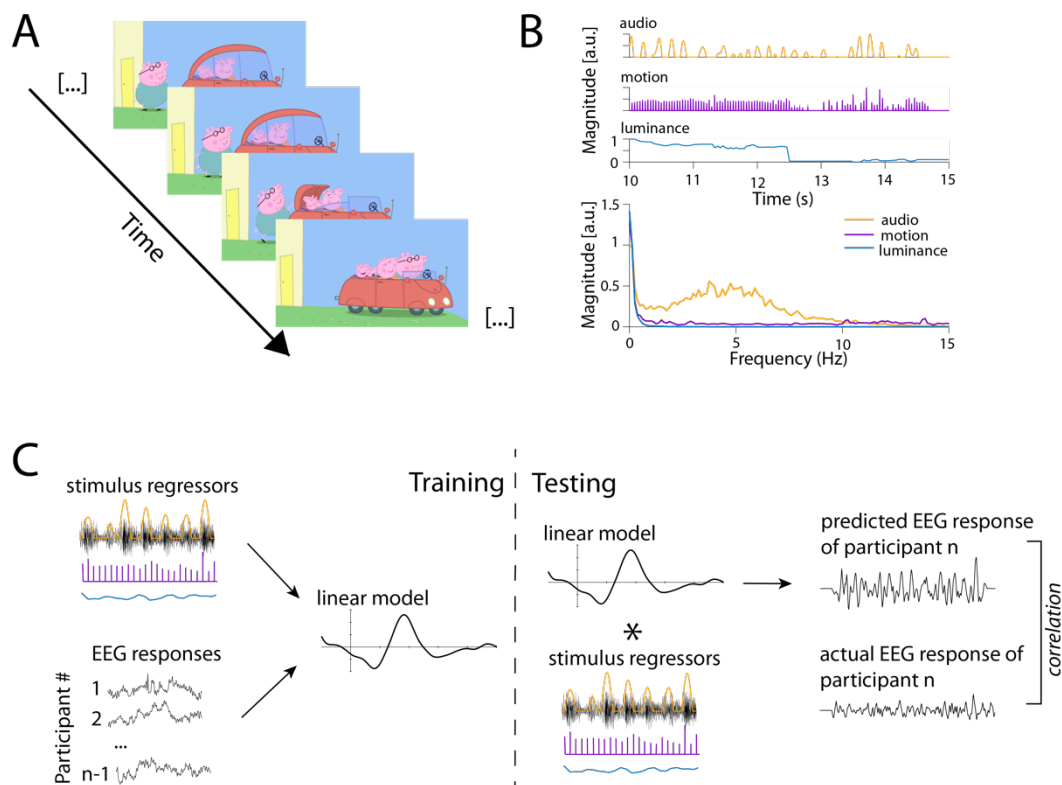
240

241    *Stimulus parameters.* As expected for a child-friendly cartoon movie, frame-to-frame

242    fluctuations in luminance were small. On average, the change in luminance from one frame

243    to the next was 0.35 units per frame (range 0–53, median = 0.05). Note that this deviates

244    from previous studies where the entire dynamic range of luminance (i.e., black to white)

245    was used to quantify the temporal response function in the adult EEG response (e.g., Lalor,

246    Pearlmutter, Reilly, McDarby, & Foxe, 2006; Vanrullen & MacDonald, 2012) or in non-

247    human animal electrophysiological responses (Ringach & Shapley, 2004). In contrast, the

248    luminance-derived motion regressor yielded sizable variance, with a mean frame-to-frame

249    change of 38 units (range 0–192, median = 36).

250    In sum, while variance in the luminance regressor was small, both motion and

251    audio regressors showed considerable and promising degrees of variance.

252    Lastly, all regressors were downsampled (audio) and interpolated (motion,

253    luminance), respectively, to the EEG sampling frequency of 250 Hz. In all regressors, time

254    periods in which no EEG data was available as a result of artefact rejection during

255    preprocessing were zero-replaced. Finally, EEG data and physical regressors were aligned

256    and available for the linear model analysis.

257



259    *Figure 1. Physical properties of stimulus regressors.* A) shows four exemplary stills from
260    the movie used as stimulus material. B) shows an example of a 5-s-long stretch of the audio
261    (orange), motion (purple), and luminance regressors (blue). Below, the frequency spectrum
262    of the stimulus regressors is depicted; while frequencies < 10 Hz appear to be dominant in
263    the audio regressors, no such dominance can be observed for the other regressors. C) shows
264    an overview of the analysis approach. During training, stimulus regressors and the EEG
265    signal of n-1 participants was used to compute a generic response function (left part).
266    During testing, this generic response function was used to predict the EEG response of the
267    nth participant, which was then compared to the actual EEG response of that participant.
268    See main text for further details.

269

270    *Temporal response functions (TRF).* To quantify the degree to which the measures EEG of

271    7-month-olds (as well as adults) can be expressed as a linear response to stimulus features,

9

272  we used regularized regression (with ridge parameter λ) as implemented in the mTRF
273  toolbox (Crosse et al., 2016). The key idea here is to estimate a temporal response function
274  (TRF), that is, a set of time-lagged weights *g*, with which a regressor *s* (here, the physical
275  stimulus features) would need to be convolved (i.e. multiplied and summed) in order to
276  optimally predict the measured EEG response *r*.

277        More specifically, we used a forward encoding model approach. In a first pass, we
278  aimed to maximize the predictive accuracy of such a model by estimating so-called
279  "generic" models, that is, we predicted the EEG data of an *n*th participant based on a
280  "generic" temporal response function (TRF) from *n-1* participants to the auditory or visual
281  stimulus signal. Since changes in the EEG signal are not likely to occur simultaneously
282  with changes in the stimulus signal but rather with an (unknown) time lag, predictions were
283  computed over a range of time lags between 200 ms earlier than the stimulus signal and
284  1000 ms later than the stimulus signal.

285  *Choosing the optimal regularization parameter λ.* To obtain the optimal regularization
286  parameter λ for each stimulus regressor separately, as well motion and audio
287  simultaneously, we trained the respective model on a variety of λ values between $10^{-5}$ and
288  $10^5$, increasing the exponent in steps of 0.5, and used the resulting models to predict the
289  EEG signal for each participant. We then computed the mean response function across *n-1*
290  participants and used this response function to predict EEG response of the nth participant
291  (i.e., *n*-fold leave-one-out crossvalidation). Finally, we computed the predictive accuracy
292  (i.e., Pearson's correlation coefficient *r* between the predicted EEG response and the actual
293  EEG response) for each participant, resulting in one accuracy value for each electrode (27
294  for infants, 64 for adults) per participant and stimulus parameter for each λ value. For each
295  participant, stimulus parameter, and electrode, we selected the λ value maximizing
296  *predictive* accuracy. Based on these values, we obtained the mean regularization parameter
297  λ value across all electrodes and participants (see Table S1).

298        These optimal λ parameters were used in the following to train the model, resulting
299  in separate response functions for each stimulus parameter. For each of the three physical
300  stimulus parameters (luminance, motion, audio) we computed a separate model. In
301  addition, we computed a model using both, motion and audio, as regressors ("joint audio-
302  motion model"). We chose not to include luminance in this model, as the regressor for
303  luminance did not yield any reliable model in itself (see results).

304  *Evaluation of temporal response functions.* For statistical evaluation of the resulting
305  response functions, we computed a cluster-based permutation test with 1000

10

randomizations, testing the obtained response functions against zero. A cluster was defined along the dimensions time and electrode position, with the constraint that a cluster had to extend over at least two adjacent electrodes. A type-1-error probability of less than .05 was ensured for all clusters.

In addition, to assess internal validity of our model predictions on an individual basis, we computed three different predictive accuracies per participant. First, for each participant *n*, we computed the correlation between the predicted response generated on a model trained on *n–1* participants and the actual EEG response of *n* ("generic model").

Second, rather than relying on the generic model based on *n–1* participants, we computed an individual response function for each participant ("individual model"). To that end, 80 % of the available data for a given participant were used to train the model, and the resulting response function was then correlated with the response observed in the remaining 20 % of the data.

Third, a *permuted* or null predictive accuracy ("shifted control") was obtained. Before calculating accuracy this way, we shifted the actual EEG response for participant *n* in steps of 2 s (in order to ensure to exceed the potential autoregressive structure of the EEG data) and computed the correlation between the shifted EEG signal and the predicted response, based on the generic model trained on *n–1* participants.
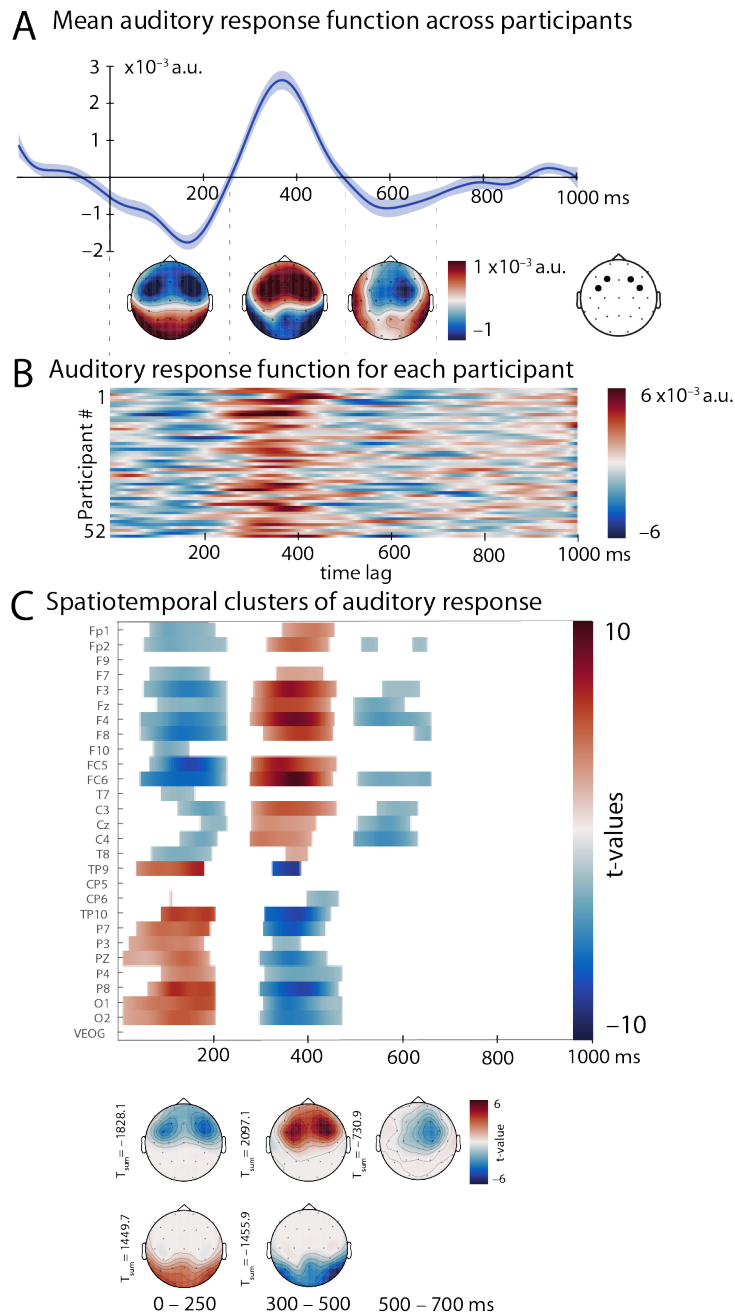
RESULTS

*Temporal response function.* We computed a generic temporal response function for each stimulus regressor as well as the audio and motion regressor combined (joint audio-motion model).

We observed a clearly defined response function using the audio regressor (Figure 2 for infants and Figure S2 for adults) and the motion regressor (Figure 3 and Figure S3 for adults), while no clear response function could be obtained using the luminance regressor for either infants or adults (Figure S1). While Figures 2 and 3 show the respective response functions obtained from a model which included both regressors (joint audio-motion model), comparable response functions resulted when using either of the regressors in isolation.

Interestingly, the observed response function did not only become visible in the average response function, but also for the vast majority of participants on an individual level (Figure 2B and 3B). Furthermore, note that while a clearly defined response was

11

339  visible for both, the audio and motion regressor, the amplitude of the response function for

340  the motion regressor was much smaller compared to the amplitude of the audio response

341  function.

342



*Figure 2. Auditory response function (using motion and audio regressor simultaneously) for infant participants.* A) shows the mean mTRF (mean ± within-subject SEM) computed across all participants, averaged over FC5, FC6, F3, and F4, and topographic representations for 0–250 ms, 250–500 ms, and 500–700 ms with electrodes included in the above-shown average marked by black dots. B) shows the auditory response function for each individual infant. C) displays the results of the cluster-based permutation test,
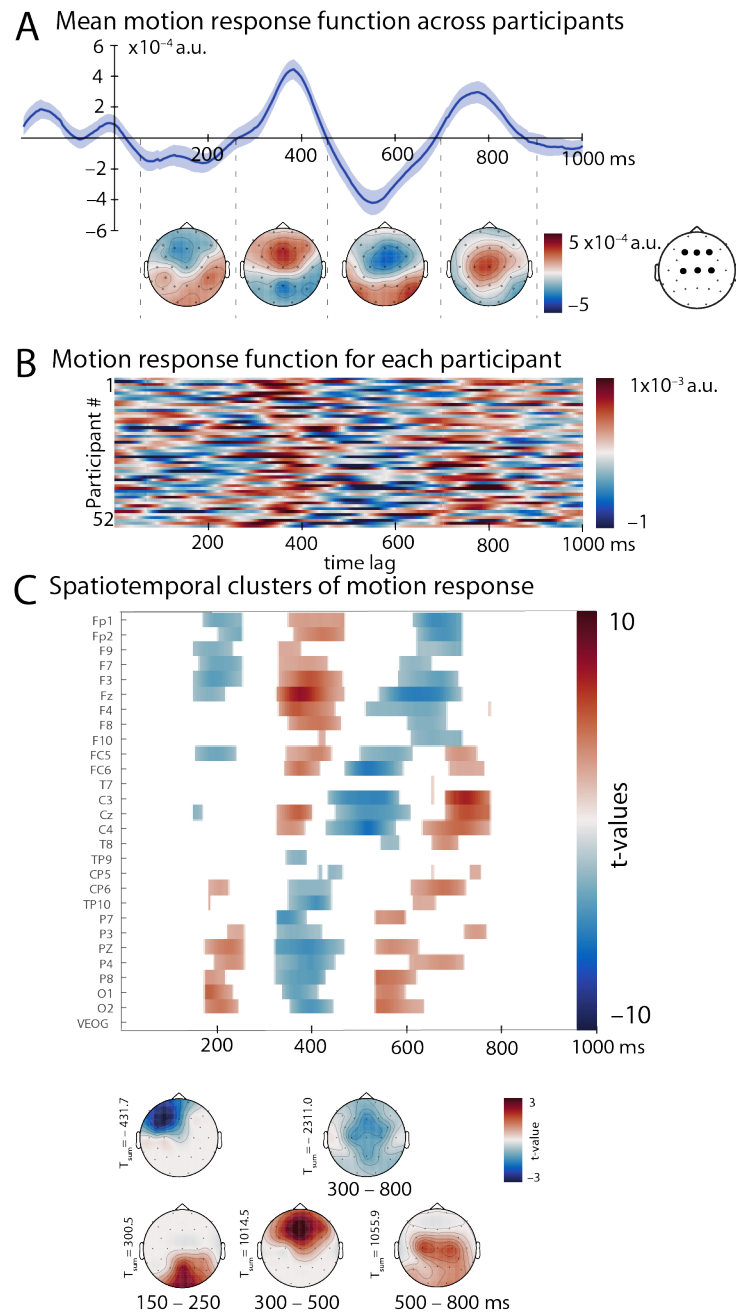
12

350 comparing the response function shown in A) and B) to zero. Positive deviations are
351 displayed in red, while negative deviations are shown in blue. In the bottom part of C), the
352 same clusters as in the top part of C) are shown as topographic distributions, along with the
353 summed t-value across the cluster.

354 *Comparing infant and adult brain responses.* When directly comparing infant and adult
355 response functions (Figure 4), similarities as well as striking differences emerge. Overall,
356 amplitudes of the response functions are comparable for infants and adults, both showing
357 the already mentioned larger amplitudes for audio regressors and smaller amplitudes for
358 motion regressors. For both, infants and adults, the auditory response function is marked
359 by a prominent frontocentral positivity (250–500 ms for infants, 300–450 ms for adults).
360 While this response appears to be slightly longer for infants, overall, both latency and
361 topography indicate a comparable response for infants and adults. In contrast, the infant
362 auditory response function lacks a second, earlier and more central positivity, which can
363 be observed between 150 and 250 ms in the adult auditory response function.

364 For the motion response function, both infants and adults show two frontal /
365 frontocentral positivities (250– 450 and 700–900 ms for infants and 250–350 and 450–550
366 ms for adults). Hence, infants and adults show a comparable response, though the infant
367 response appears to be much slower and less temporally modulated.

368 *Cluster-based permutation test.* We computed a cluster-based permutation test comparing
369 the temporal response function obtained using the motion, luminance, and audio regressor
370 as well as the motion and audio regressor simultaneously. We did not observe any
371 significant cluster using the luminance regressor for either infants or adults. In contrast, we
372 did obtain multiple significant clusters, indicating a positive or negative deviation from
373 zero, for the motion and audio regressor, both when included separately as well as in
374 combination (see supplementary material for a full list of the results of the cluster-based
375 permutation test using audio and motion regressor separately as well as in combination for
376 infants and adults, Figure 2C and 3C for infant results and S2C and S3C for adult results).
377 The resulting clusters confirm the deflections observed in the auditory and motion response
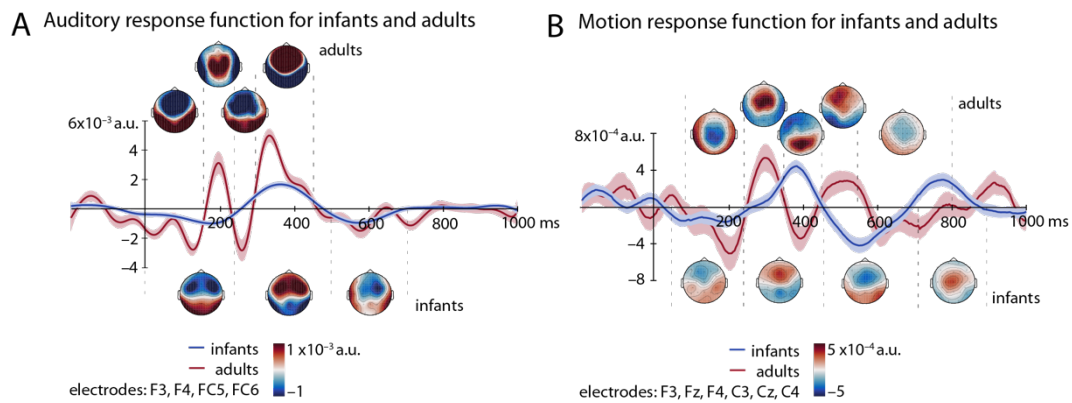378 function (Figure 2A and 3A, respectively).

379

13

*Figure 3. Motion response function (using motion and audio regressor simultaneously) for infant participants.* A) shows the mean mTRF (mean ± SEM) computed across all participants, averaged over F3, Fz, F4, C3, Cz, and C4, and topographic representations for 50–250 ms, 250–450 ms, 450–700 ms, and 700–900 ms with electrodes included in the above-shown average marked by black dots. B) shows the motion response function for each individual infant. C) displays the results of the cluster-based permutation test, comparing the response function shown in A) and B) to zero. Positive deviations are displayed in red, while negative deviations are shown in blue. In the bottom part of C), the same clusters as in the top part of C) are shown as topographic distributions, along with the summed t-value across the cluster.

*Figure 4. Comparison of infant and adult response functions.* Mean mTRF for infant (in blue) and adult (in red) participants are shown for the audio regressor (A) and the motion regressor (B). The infant response functions and topographical representations are identical to those shown in Fig 2A and 3A for audio and motion regressors, respectively. Responses are averaged across the same electrodes for adults and infants, namely FC5, FC6, F3, and F4 for A) and F3, Fz, F4, C3, Cz, and C4 for B). The topographic representations of adult responses correspond to those in the supplementary material, namely 50–150 ms, 150–250 ms, 250–300 ms, and 300–450 ms for A) and 50–250 ms, 250–350 ms, 350–450 ms, 450–550 ms, and 550–800 ms for B).
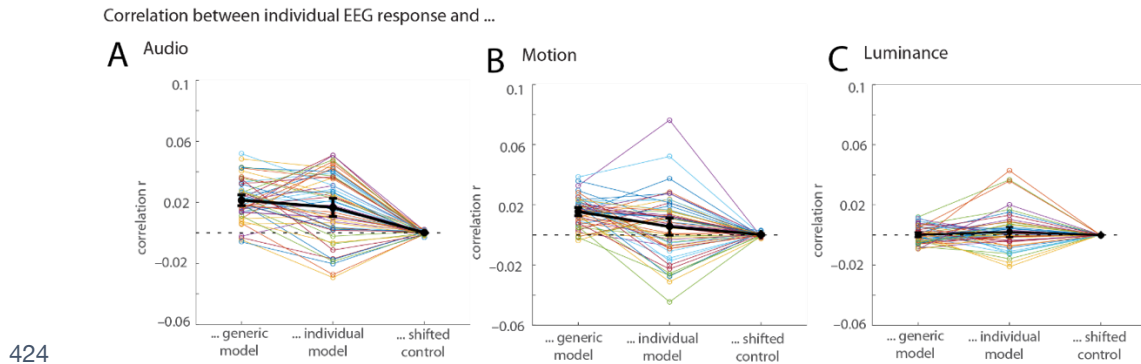
*Generic vs. individual response functions.* The results discussed above rely on a generic model computed based on data from *n–1* participants in order to predict the *n*th participant (see Di Liberto & Lalor, 2017). An alternative approach (and in fact preferable, if enough data for per subject is available; e.g., Fiedler et al., 2019; O'Sullivan et al., 2017) computes an individual model based on a subset of an individual's data and compare the resulting predictions to the remaining data.

As expected, individual models showed a larger variance compared to the generic model (Figure 5–7; see S4 and S5 for data from adult participants), but both, generic model and individual model result in correlations clearly above zero (with the exception of luminance, where no reliable prediction was possible for either mode, see Figure 5C).

When both, audio and motion regressor were included (Figure 6), the generic model resulted in a higher correlation compared to the individual model for infant participants ($t(51)=3.76$, $p<.001$); 37 participants showed a higher correlation with the generic model while only 15 participants showed a higher correlation with the individual model. When using only the motion regressor (Figure 5B), the correlations were also higher for the generic compared to the individual model ($t(51)=3.50$, $p<.001$), while for the auditory regressor (Figure 5A), this difference was less pronounced ($t(51)=1.82$, $p=.07$).
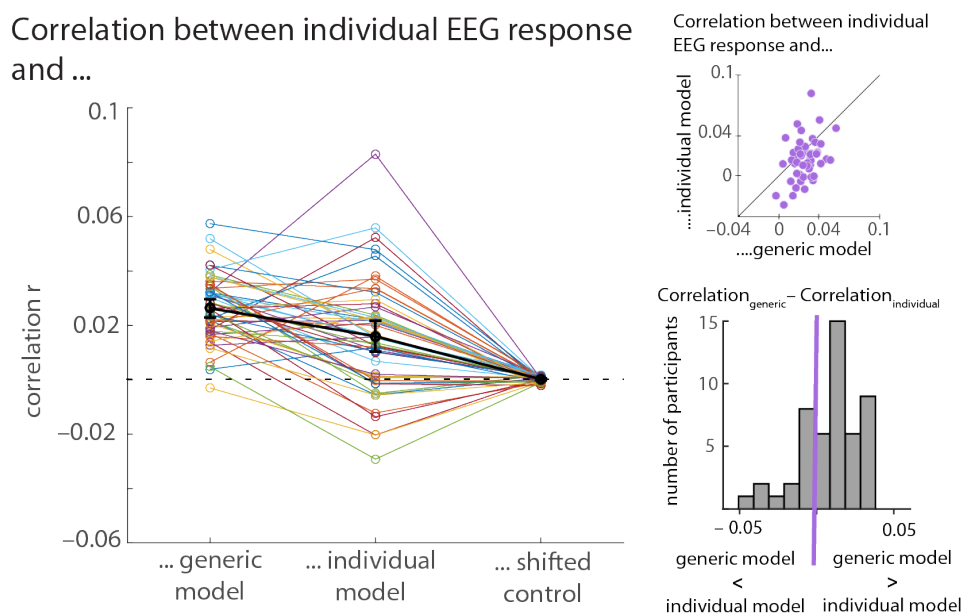
15

420     As a control analysis, a generic model using temporally shifted (i.e., purposefully

421     misaligned) versions of the actual EEG signal (1,000 iterations) did yield substantially

422     lower predictive accuracy values.

423



*Figure 5. Predictive Accuracy (r) between model and EEG response for infant participants.*
The recorded individual EEG response was correlated with three different parameters using
Pearson's correlation coefficient for the audio regressor (A), motion regressor (B), and
luminance regressor (C). On the left, the correlation between the recorded EEG responses
of participant n and the response predicted by the generic model based on the remaining n-
1 participants is shown for each participant. In the middle, the correlation between the
model trained on the first 80 % of the data available for each participant and used to predict
the remaining 20 % from that participant and the actual EEG response recoded from that
participant is shown. The right column shows the correlation between the prediction
generated by the generic model and the recorded EEG data shifted in a circular way in steps
of 2 s as a control condition (averaged over all possible shifts). Correlations are shown for
each infant participant (in colors) as well as the mean correlation with 95% CI (confidence
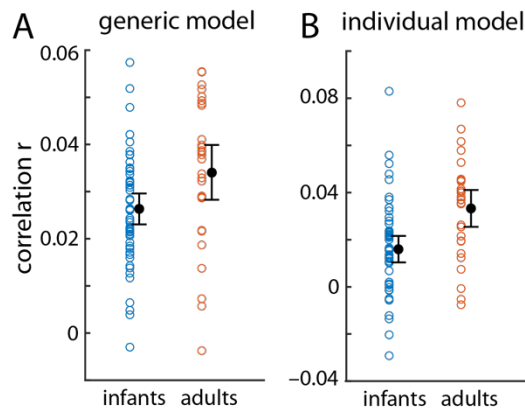interval) across all participants (in black).

438



16

440 *Figure 6. Predictive accuracy for modelled and observed EEG response for infant*
441 *participants in a joint audio–motion model.* The left part of the figure shows the correlation
442 (based on Pearson's correlation coefficient) between the recorded EEG signal and the EEG
443 responses predicted based on the generic model (left column), the individual model (middle
444 column), and a shifted control condition (right column, see text). The two plots on the right
445 hand visualize a comparison between the generic and the individual model. In the top plot,
446 each purple dot indicates the difference between the correlation with the generic model and
447 the correlation with the individual model. Hence, a purple dot in the right bottom part of
448 the graph indicates an individual with a higher correlation for the generic compared to the
449 individual model, while a purple dot in the top left part indicates an individual with a higher
450 correlation for the individual compared to the generic model. The bottom plot displays the
451 same information in a bar graph; individuals having a higher correlation for the generic
452 model have a positive difference and hence fall to the right of the zero-threshold marked in
453 purple while those with a higher correlation for the individual model have a negative
454 difference and fall to the left of the zero-threshold.

455



456

457 *Figure 8. Predictive accuracy for infants and adults in a joint audio–motion model.* A)
458 shows the individual correlations using Pearson's correlation coefficient for infants (blue)
459 and adults (orange) using the generic model. B) shows the individual correlations using
460 Pearson's correlation coefficient for infants (blue) and adults (orange) using the individual
461 model. Mean accuracies with 95 % confidence intervals are shown in black.

462

463 DISCUSSION

464 We investigated the use of a variant of forward encoding models (multivariate temporal

465 response functions, mTRFs) to analyze infant brain responses to a continuous complex

466 audiovisual stimulus, namely a 5-minute cartoon movie. We observed clearly defined

467 response patterns to both the auditory as well as the motion content, but no predictive

468 response function for changes in luminance was found.

469      Our results demonstrate that the simultaneous acquisition of individual brain

470 responses to different sensory modalities is possible in the infant brain, opening new

471 avenues for ecologically valid multisensory research paradigms in developmental

472 neuroscience. Furthermore, our results suggest that a generic model derived from a larger

17

473    set of unrelated infant data is as good or slightly better compared to an individual model in

474    predicting the individual brain response, especially in cases where only limited data is

475    available. This points to the further utility of such an approach in developmental and at-

476    risk populations.

477    *Motion and Audio.* For both, motion and audio information in the cartoon movie, we found

478    a clearly defined response in both infants and adults. The observed responses are largely

479    consistent with patterns typically reported in more traditional event-related brain potentials.

480    The frontocentral negativity between 450 and 700 ms for instance observed in the infant

481    brain responses linked to the motion regressor corresponds in timing, shape, and

482    topography to the Nc component, an infant ERP component that can routinely be observed

483    in visual paradigms and has been linked to attention allocation (Webb et al., 2005).

484    Likewise, the bifocal frontal positivity observed in the infants' brain response linked to the

485    auditory envelope shows a strong similarity to the commonly reported P2 response in infant

486    auditory brain responses (Wunderlich et al., 2006).

487    The direct comparison of infant and adult brain responses (Figure 4) may provide

488    insight into the developmental changes. In response to the auditory envelope, both infants

489    and adults show a prominent frontal negativity peaking around 400 ms. Notably, however,

490    the adults show an additional central positivity around 200 ms, which is missing in the

491    infant response. This corroborates and replicates known developmental changes commonly

492    observed in auditory evoked responses when comparing infants and adults (Wunderlich &

493    Cone-Wesson, 2006). Considering the motion response, the correspondence between infant

494    and adult response is less straight-forward. While the adult response is characterized by

495    two frontocentral positivities, one peaking around 300 ms and the other around 500 ms, the

496    infant response is dominated by one frontocentral peak around 400 ms.

497    Importantly, we used both, generic response functions as well as individual

498    response functions to predict the EEG signal. When using both, the motion and the auditory

499    regressor, performance was significantly better for the generic compared to the individual

500    model. When using only the auditory regressor, the same pattern was visible but the

501    difference only marginally significant. Note, however, that both, generic and individual

502    models generated predictions that were significantly above chance level. This demonstrates

503    two important things. First, five minutes of EEG recording are sufficient to compute

504    reliable models, both on an individual level as well as across participants as a generic

505    model. This is not only true for EEG data obtained from healthy adults but also for data

506    obtained from populations providing notoriously noisy signal, such as infants. Second,

507    brain responses across participants, both infants and adults, are sufficiently similar to

18

508   generate a model that can successfully predict a new infant's brain response, yielding even

509   better outcomes compared to the individual model.

510   *Limitations and future studies*. The present study provides an important step and proof of

511   feasibility for using mTRFs to analyze infant EEG data in response to complex and

512   dynamic audiovisual stimulus material. This offers a whole host of new possibilities in the

513   investigation of infant's brain responses in their natural environment.

514   One important feature of the present study is that we used the unmanipulated

515   cartoon video material. While this makes for an ecologically valid and easy-to-obtain

516   stimulus, it comes with the caveat of a lack of control for stimulus properties.

517   Notably, while we did observe a clear-cut response to the motion and the auditory

518   regressor, we did not find a reliable response to the changes in luminance. The most likely

519   explanation for this discrepancy is the lack in variance in the luminance content. While the

520   motion and the auditory regressor showed large-amplitude changes throughout the video

521   (e.g., average motion change between frames = 38 units), average luminance of this cartoon

522   movie remained fairly constant (average luminance change between frames = 0.35 units).

523   Previous studies targeting neural responses to luminance change (in adults) typically used

524   considerably more pronounced black–white contrast (Lalor et al., 2006; Vanrullen &

525   MacDonald, 2012). Hence, the luminance changes in the stimulus material were likely too

526   small to elicit any robust change in brain response. Future studies explicitly varying the

527   luminance content are therefore necessary to investigate the applicability of mTRFs to

528   other visual stimulus parameters in infants.

529   Also, we operationalized motion as change in pixel from one frame to the next.

530   This means that the motion regressor not only reflected the actual motion of the objects and

531   persons depicted in the video but also cuts in the video. For the present purpose, we did not

532   differentiate between these two possibilities of motion.

533   Building upon the present results, a next step would therefore be to purposefully

534   manipulate such parameters. By using stimulus material designed to encompass a larger

535   variance in luminance and/or no cuts in the video, it should for instance be possible to

536   observe brain response to changes in luminance and motion responses that can be clearly

537   linked to actual motion rather than video cuts. Such an approach could for instance provide

538   valuable new insights into the processing of biological motion (Marshall & Shipley, 2009;

539   Reid, Hoehl, & Striano, 2006).

540   Furthermore, in the present study, we did not contrast different conditions, neither

541   within infant nor between different groups of infants. Having demonstrated the feasibility

19

of using encoding models to model brain responses for this type of complex audiovisual stimuli, the next step would certainly be to utilize this approach to investigate differences in processing between (a) different types of stimulation or (b) different groups of infants.

A first step in using mTRFs to contrast different continuous stimulus signals has been done by Kalashnikova et al. (2018), who compared the processing of infants vs. adult directed speech in 7-month-olds. Future studies could encompass more complex naturalistic scenarios, using for instance audiovisual video material. More importantly, mTRFs can also be used to investigate brain responses in live interactions, in which the live input the infant receives is recorded and used as a regressor in the subsequent analysis. Such an approach would provide an important tool in investigating the neural bases of social interactions.

*Conclusion*. The present data demonstrate that forward encoding models based on the multivariate temporal response function (mTRF) pose a valuable and versatile tool in quantifying and disentangling complex audiovisual brain responses and the according perceptual processes in infancy. Our results open way for applications to a variety of research areas not only in early development, but also in other special populations characterized by short attention spans and low cooperativeness, including research in severely impaired neurological patients. New paradigms could not only entail complex multisensory perception, but extend to dynamic social interactions. As such, mTRF approaches to infant data analysis will allow developmental researchers to devise more engaging and thereby more easily applicable experimental set-ups for infancy research.

DATA ACCESSIBILITY.

Data will be made available upon publication.

571    BIBLIOGRAPHY

572    Bartels, A., Zeki, S., & Logothetis, N. K. (2008). Natural vision reveals regional
573        specialization to local motion and to contrast-invariant, global flow in the human
574        brain. *Cerebral Cortex*. https://doi.org/10.1093/cercor/bhm107

575    Broderick, M. P., Anderson, A. J., Di Liberto, G. M., Crosse, M. J., & Lalor, E. C.
576        (2018). Electrophysiological Correlates of Semantic Dissimilarity Reflect the
577        Comprehension of Natural, Narrative Speech. *Current Biology*.
578        https://doi.org/10.1016/j.cub.2018.01.080

579    Crosse, M. J., Di Liberto, G. M., Bednar, A., & Lalor, E. C. (2016). The Multivariate
580        Temporal Response Function (mTRF) Toolbox: A MATLAB Toolbox for Relating
581        Neural Signals to Continuous Stimuli. *Frontiers in Human Neuroscience*.
582        https://doi.org/10.3389/fnhum.2016.00604

583    Dayan, P., & Abbott, L. (2001). *Theoretical Neuroscience: Computational and*
584        *Mathematical Modeling of Neural Systems.* Cambridge, MA: MIT Press.

585    de Haan, M., Johnson, M. H., & Halit, H. (2003). Development of face-sensitive event-
586        related potentials during infancy: a review. *International Journal of*
587        *Psychophysiology*, *51*(1), 45–58. Retrieved from
588        http://www.ncbi.nlm.nih.gov/pubmed/14629922

589    Di Liberto, G. M., & Lalor, E. C. (2017). Indexing cortical entrainment to natural speech
590        at the phonemic level: Methodological considerations for applied research. *Hearing*
591        *Research*. https://doi.org/10.1016/j.heares.2017.02.015

592    Ding, N., & Simon, J. Z. (2013). Adaptive Temporal Encoding Leads to a Background-
593        Insensitive Cortical Representation of Speech. *Journal of Neuroscience*.
594        https://doi.org/10.1523/jneurosci.5297-12.2013

595    Ellis, C. T., & Turk-Browne, N. B. (2018). Infant fMRI: A Model System for Cognitive
596        Neuroscience. *Trends in Cognitive Sciences*.
597        https://doi.org/10.1016/j.tics.2018.01.005

598    Fiedler, L., Wöstmann, M., Graversen, C., Brandmeyer, A., Lunner, T., & Obleser, J.
599        (2017). Single-channel in-ear-EEG detects the focus of auditory attention to
600        concurrent tone streams and mixed speech. *Journal of Neural Engineering*.
601        https://doi.org/10.1088/1741-2552/aa66dd

602    Fiedler, L., Wöstmann, M., Herbst, S. K., & Obleser, J. (2019). Late cortical tracking of

ignored speech facilitates neural selectivity in acoustically challenging conditions. *NeuroImage*. https://doi.org/10.1016/j.neuroimage.2018.10.057

Hamilton, L. S., & Huth, A. G. (2018). The revolution will not be controlled: natural stimuli in speech neuroscience. *Language, Cognition and Neuroscience*. https://doi.org/10.1080/23273798.2018.1499946

Hasson, U., Nir, Y., Levy, I., Fuhrmann, G., & Malach, R. (2004). Intersubject Synchronization of Cortical Activity during Natural Vision. *Science*. https://doi.org/10.1126/science.1089506

Huk, A., Bonnen, K., & He, B. J. (2018). Beyond Trial-Based Paradigms: Continuous Behavior, Ongoing Neural Activity, and Natural Stimuli. *The Journal of Neuroscience*. https://doi.org/10.1523/jneurosci.1920-17.2018

Jessen, S., & Kotz, S. A. (2011). The temporal dynamics of processing emotions from vocal, facial, and bodily expressions. *NeuroImage*, *58*(2), 665–674.

Jones, E. J. H., Venema, K., Lowy, R., Earl, R. K., & Webb, S. J. (2015). Developmental changes in infant brain activity during naturalistic social experiences. *Developmental Psychobiology*. https://doi.org/10.1002/dev.21336

Kalashnikova, M., Peter, V., Di Liberto, G. M., Lalor, E. C., & Burnham, D. (2018). Infant-directed speech facilitates seven-month-old infants' cortical tracking of speech. *Scientific Reports*. https://doi.org/10.1038/s41598-018-32150-6

Lalor, E. C., Pearlmutter, B. A., Reilly, R. B., McDarby, G., & Foxe, J. J. (2006). The VESPA: A method for the rapid estimation of a visual evoked potential. *NeuroImage*. https://doi.org/10.1016/j.neuroimage.2006.05.054

Leong, V., Byrne, E., Clackson, K., Georgieva, S., Lam, S., & Wass, S. (2017). Speaker gaze increases information coupling between infant and adult brains. *Proceedings of the National Academy of Sciences*. https://doi.org/10.1073/pnas.1702493114

Marshall, P. J., & Shipley, T. F. (2009). Event-related potentials to point-light displays of human actions in 5-month-old infants. *Developmental Neuropsychology*. https://doi.org/10.1080/87565640902801866

Matusz, P. J., Dikker, S., Huth, A. G., & Perrodin, C. (2018). Are we ready for real-world neuroscience? *Journal of Cognitive Neuroscience*. https://doi.org/10.1162/jocn_e_01276

634  Naselaris, T., Kay, K. N., Nishimoto, S., & Gallant, J. L. (2011). Encoding and decoding
635      in fMRI. *NeuroImage*. https://doi.org/10.1016/j.neuroimage.2010.07.073

636  Nishimoto, S., Vu, A. T., Naselaris, T., Benjamini, Y., Yu, B., & Gallant, J. L. (2011).
637      Reconstructing visual experiences from brain activity evoked by natural movies.
638      *Current Biology*. https://doi.org/10.1016/j.cub.2011.08.031

639  O'Sullivan, A. E., Crosse, M. J., Di Liberto, G. M., & Lalor, E. C. (2017). Visual Cortical
640      Entrainment to Motion and Categorical Speech Features during Silent Lipreading.
641      *Frontiers in Human Neuroscience*. https://doi.org/10.3389/fnhum.2016.00679

642  Oostenveld, R., Fries, P., Maris, E., & Schoffelen, J.-M. (2011). FieldTrip: Open source
643      software for advanced analysis of MEG, EEG, and invasive electrophysiological
644      data. *Computational Intelligence and Neuroscience*, *2011*, 156869.

645  Pichon, S., de Gelder, B., & Grèzes, J. (2009). Two different faces of threat. Comparing
646      the neural systems for recognizing fear and anger in dynamic body expressions.
647      *NeuroImage*, *47*, 1873–1883.

648  Reid, V. M., Hoehl, S., & Striano, T. (2006). The perception of biological motion by
649      infants: An event-related potential study. *Neuroscience Letters*.
650      https://doi.org/10.1016/j.neulet.2005.10.080

651  Reynolds, G. D., & Guy, M. W. (2012). Brain-behavior relations in infancy: Integrative
652      approaches to examining infant looking behavior and event-related potentials.
653      *Developmental Neuropsychology*. https://doi.org/10.1080/87565641.2011.629703

654  Ringach, D., & Shapley, R. (2004). Reverse correlation in neurophysiology. *Cognitive*
655      *Science*. https://doi.org/10.1016/j.cogsci.2003.11.003

656  Ru, P. (2001). *Multiscale Multirate Spectro-Temporal Auditory Model*. University of
657      Maryland College Park.

658  Stets, M., Stahl, D., & Reid, V. M. (2012). A meta-analysis investigating factors
659      underlying attrition rates in infant ERP studies. *Developmental Neuropsychology*,
660      *37*(3), 226–252. https://doi.org/10.1080/87565641.2012.654867

661  Vanrullen, R., & MacDonald, J. S. P. (2012). Perceptual echoes at 10 Hz in the human
662      brain. *Current Biology*. https://doi.org/10.1016/j.cub.2012.03.050

663  Webb, S. J., Long, J. D., & Nelson, C. A. (2005). A longitudinal investigation of visual
664      event-related potentials in the first year of life. *Dev Sci*, *8*(6), 605–616.

665        https://doi.org/10.1111/j.1467-7687.2005.00452.x

666    Wunderlich, J. L., & Cone-Wesson, B. K. (2006). Maturation of CAEP in infants and

667        children: A review. *Hearing Research*. https://doi.org/10.1016/j.heares.2005.11.008

668    Wunderlich, J. L., Cone-Wesson, B. K., & Shepherd, R. (2006). Maturation of the

669        cortical auditory evoked potential in infants and young children. *Hearing Research*.

670        https://doi.org/10.1016/j.heares.2005.11.010

671