

1 **Cortical tracking of visual rhythmic speech by 5- and 8-month-old infants: Individual**
2 **differences in phase angle relate to language outcomes up to 2 years**

3

4 Áine Ní Choisdealbha, Adam Attaheri, Sinead Rocha, Natasha Mead, Helen Olawole-Scott,
5 Maria Alfaro e Oliveira, Carmel Brough, Perrine Brusini, Samuel Gibbon, Panagiotis Boutris,
6 Christina Grey, Isabel Williams, Sheila Flanagan, and Usha Goswami *

7

8 Centre for Neuroscience in Education, Department of Psychology, University of Cambridge

9

10 *Corresponding author: Usha Goswami, ucg10@cam.ac.uk

11

12 **Running Title: *Visual Phase Angle and Language Outcomes***

13

14 **Acknowledgements:**

15 We are grateful to the families in the Cambridge BabyRhythm UK project for their
16 participation. This work was funded by the European Research Council (ERC) under the
17 European Union's Horizon 2020 research and innovation programme (grant agreement
18 No.694786).

19 Data availability statement: Scripts and processed data used for analysis will be made
20 available at <https://osf.io/tmuqy/> upon acceptance

21 This work was reviewed by the Psychology Research Ethics Committee of the University of
22 Cambridge.

23 The authors declare that they have no competing interests.

24

25

26

27 **Highlights**

- 28 • Infant preferred phase to visual rhythmic speech predicts language outcomes
- 29 • Significant cortical tracking of visual speech is present at 5 and 8 months
- 30 • Phase angle to visual speech at 8 months predicted greater receptive and productive
- 31 vocabulary at 24 months.

32 **Abstract**

33 It is known that the rhythms of speech are visible on the face, accurately mirroring changes in
34 the vocal tract. These low-frequency visual temporal movements are tightly correlated with
35 speech output, and both visual speech (for example, mouth motion) and the acoustic speech
36 amplitude envelope entrain neural oscillations. Low-frequency visual temporal information
37 (‘visual prosody’) is known from behavioural studies to be perceived by infants, but oscillatory
38 studies are currently lacking. Here we measure cortical tracking of low-frequency visual
39 temporal information by five- and eight-month-old infants using a rhythmic speech paradigm
40 (repetition of the syllable “ta” at 2 Hz). Eye-tracking data was collected simultaneously with
41 EEG, enabling computation of cortical tracking and phase angle during visual-only speech
42 presentation. Significantly higher power at the stimulus frequency indicated that cortical
43 tracking occurred across both ages. Further, individual differences in preferred phase to visual
44 speech related to subsequent measures of language acquisition. The difference in phase
45 between visual-only speech and the same speech presented as auditory-visual at 6- and 9-
46 months was also examined. These neural data suggest that individual differences in early
47 language acquisition may be related to the phase of entrainment to visual rhythmic input in
48 infancy.

49 **Keywords:** infant, EEG, visual speech, language acquisition

50

51 **1. Introduction**

52 Visual speech inputs are known to enhance auditory speech information for adults, particularly
53 in noisy or degraded conditions (e.g., Sumbly & Pollack, 1954; Grant & Braida, 1991; Grant &
54 Seitz, 2000). Visual speech information may also enhance speech processing for infants, as
55 already demonstrated behaviourally regarding phonetic information, which can be easy to see
56 on the face (Kuhl & Meltzoff, 1982; Pons et al., 2009; Walden et al., 1977). The theoretical
57 focus adopted here concerns speech prosody and speech rhythm rather than phonetic
58 information, as the acoustic statistics of infant-directed speech (IDS) foreground speech rhythm
59 patterns (Leong et al., 2017). Speech rhythm is also visible on the face (Munhall, Kross &
60 Vatikiotis-Bateson, 2002; Kitamura, Guellai & Kim, 2014). Indeed, the temporal
61 correspondence between visual and auditory prosody is highlighted by the natural statistics of
62 audiovisual (AV) speech, as lip, jaw, cheek and head movements convey crucial information
63 about the speech amplitude envelope (AE), which carries acoustic prosodic information
64 (Chandrasekaran, Trubanova, Stillitano, Caplier & Ghazanfar, 2009). Low-frequency spatial
65 and temporal modulations in the 2 – 7 Hz range related to mouth opening and closing are of
66 particular importance for AE tracking, causing Chandrasekeran et al. (2009) to observe that the
67 natural temporal features of AV speech signals are optimally structured for the neural
68 (oscillatory) rhythms of the brains of their receivers. In the current report, we investigate
69 whether these oscillatory rhythms are encoding visual prosody in the pre-verbal infant brain,
70 and whether individual differences in the tracking of visual speech have implications for later
71 language acquisition. Speech prosody and speech rhythm play central roles in infant language
72 acquisition (Mehler et al., 1988), but until recently similarities between neural speech
73 processing mechanisms in infants and adults have been relatively under-explored (Attaheri et
74 al., 2022a,b; Tan et al., 2022).

75 There are many adult studies demonstrating oscillatory cortical tracking of the speech
76 AE for both auditory-only and AV speech (see Giraud & Poeppel, 2012, for review). Consistent
77 with the view that *visual speech* plays an important role in this AE tracking, Bourguignon and
78 colleagues (2020) showed that for adults listening passively to AV natural speech, neural
79 entrainment to the speech AE and to mouth opening/lip movements was tightly coupled (see
80 also Park et al., 2016). Further, when visual-only speech was presented (via silent videos), there
81 were significant low-frequency neural responses in auditory cortex, similar to the responses
82 recorded when speech was heard as well as seen. Bourguignon and colleagues (2020)
83 concluded that auditory neural oscillatory entrainment to prosodic acoustic information in the
84 speech AE occurs even when viewing visual-only speech information. Recent oscillatory
85 studies of speech processing by infants have demonstrated robust low-frequency cortical
86 tracking of the speech AE from birth throughout the first year (Kalashnikova et al., 2019; Ortiz
87 Barajas et al., 2021; Attaheri et al., 2022a; Tan et al., 2022), in response to both auditory-only
88 and AV inputs. Whether robust cortical tracking occurs during visual-only speech processing
89 by infants of multiple ages is currently unknown. Many of the neural mechanisms for speech
90 processing first reported in the adult literature appear to be recruited by the infant brain (e.g.,
91 theta-gamma phase amplitude coupling, Attaheri et al., 2022b). Accordingly, other
92 mechanisms found in adults, such as entrainment to visual speech, may be an intrinsic part of
93 the neural speech processing architecture. Visual articulatory cues can precede vocalisations
94 by 200ms or more, and visual gestures such as mouth opening may reset auditory cortex to the
95 “optimal state” for processing the succeeding vocalisations (Schroeder et al., 2008). Such
96 mechanisms might thus be part of the infant toolkit for language learning, and these automatic
97 neural mechanisms may help to scaffold language acquisition.

98 To investigate this possibility, here we chose to examine infant neural processing of
99 visual-only speech using a rhythmic speech paradigm. The current study is part of the

100 Cambridge UK BabyRhythm study, which tests key tenets of temporal sampling (TS) theory
101 (Goswami, 2011; 2019) in infants using a range of rhythmic measures. TS theory was originally
102 derived from neural studies of children with developmental dyslexia, who present primarily
103 with phonological (speech-sound) processing difficulties. TS theory focused first on the neural
104 processing of speech rhythm in dyslexia, identifying delta-band differences in the phase
105 alignment of low frequency neural oscillations to speech between children with and without
106 developmental dyslexia in rhythmic syllable repetition tasks (Power et al., 2013). At the
107 individual level, the angle of phase alignment in these tasks related to phonological processing
108 and reading ability (Power et al., 2012, 2013). Related dyslexia studies replicated this atypical
109 phase effect for rhythmic syllable repetition (Keshavarzi et al., 2022a), and also indicated less
110 accurate encoding of delta-band speech envelope information in connected speech tasks using
111 an mTRF approach (Power et al., 2016; Keshavarzi et al., 2022b). The Cambridge UK
112 BabyRhythm study thus set out to use similar tasks and measures with typically-developing
113 infants, and has reported on cortical tracking of sung speech using the mTRF (Attaheri et al.,
114 2022a,b) and on phase relations with rhythmic speech using a syllable repetition task in
115 auditory-only and auditory-visual (AV) modalities (Ní Choisdealbha et al., 2022, 2023).
116 Although visual-only rhythmic speech was studied in children with dyslexia, no phase-related
117 differences were found in comparison to typically-developing children (Power et al., 2013).
118 Nevertheless, it is possible that individual differences in phase alignment to visual-only
119 rhythmic speech might affect the early development of language, before infants begin to speak,
120 and we investigate that possibility here.

121 There is one prior study of infant cortical tracking of visual-only speech, using IDS and
122 the mTRF analysis approach (Tan et al., 2022). Tan and colleagues reported no significant
123 tracking of visual-only IDS by infants aged five months, despite finding significant tracking of
124 auditory-only and AV IDS by the same infants. However, Tan et al. reported that infants who

125 paid more attention to the *mouth area* during visual-only speech also showed more accurate
126 cortical tracking of the stimulus. This appears to suggest that greater attention to the rhythmic
127 cue of mouth opening facilitates tracking of visual speech by infants, which could indicate the
128 presence of cortical tracking in some participants. It is known that rhythmic facial and head
129 movements convey prosodic information about the speech AE, because head, eyebrow, lip, jaw
130 and cheek movements are systematically related to speech amplitude and fundamental
131 frequency (e.g. Munhall et al., 2002; Munhall et al., 2004). When a person is speaking, the
132 vocal tract is the source of facial movement, thus speech acoustics can be estimated from face
133 motion (and vice versa) with high reliability (Yehia, Rubin & Vatikiotis-Bateson, 1998; Yehia,
134 Kuratate & Vatikiotis-Bateson, 2002). As information about the AE is visible on the face of a
135 speaker, this information should also be available to the pre-verbal brain, making the negative
136 result reported by Tan et al. (2022) surprising. The rhythmic speech paradigm employed here
137 is a simpler input, and may be able to detect significant cortical tracking of visual-only speech
138 even at 5 months of age.

139 The rhythmic speech paradigm developed for our infants was modelled on the original
140 work by Power et al. (2012), involving repetition of the syllable “ba” at a rate of 500 ms (2
141 Hz). Note that a consideration of frequency rates is critical regarding cross-modal speech
142 dynamics, as most visual speech information is recovered from the lower spatial and temporal
143 frequencies (below 7 cycles of visual resolution per face, and 6 – 9 frames per second,
144 respectively), and adult studies of the natural statistics of AV speech highlight the importance
145 of low-frequency rhythms between 2 and 7 Hz (Chandrasekaran et al., 2009). For example,
146 when Chandrasekaran and colleagues studied natural spoken sentences in English and French,
147 they analysed the area of mouth opening as a function of time (visual temporal content). They
148 reported a close temporal correspondence between the area of mouth opening and the wideband
149 acoustic envelope, with a distinct rhythm that was between 2 and 7 Hz. Our prior neural data

150 with the same infants studied here but using AV sung speech has also highlighted the
151 importance of low-frequency modulations in the AE for cortical tracking (Attaheri et al.,
152 2022a). Infants aged four, seven and eleven months show enhanced cortical tracking in the
153 delta (0.5 – 4 Hz) and theta (4 – 8 Hz) bands compared to higher frequencies (8 – 12 Hz). Low-
154 frequency tracking of visual-only speech has yet to be explored longitudinally, and we provide
155 such information here.

156 If the human auditory speech processing architecture is engaged automatically by visual
157 speech inputs, and given our prior child data with the rhythmic speech paradigm, then pre-
158 verbal infants should also benefit from visual speech when processing human vocalisations.
159 For example, the phase alignment (or phase angle) between a stimulus and the oscillatory
160 neural response to that stimulus may relate to later language acquisition. Adult studies using
161 natural speech suggest that one key mechanism regarding efficient speech processing is the
162 *phase dynamics* of auditory and visual low-frequency neuronal oscillations (Luo, Liu and
163 Poeppel, 2010). In their MEG experiment, when adult participants were presented with
164 congruent auditory and visual information, the *visual* stream modulated *auditory* cortical
165 activity by aligning auditory oscillatory phase to the *optimal phase angle*, so that the expected
166 auditory input arrived during a high excitability state. Whether the brains of pre-verbal infants
167 would show similar cross-modal neural alignment mechanisms is not currently known.
168 However, earlier data from our BabyRhythm cohort has shown that phase alignment of neural
169 oscillations to auditory speech input is indeed active during infant speech processing, with
170 notable individual differences. For auditory-only rhythmic speech (the syllable repetition task),
171 the phase of infants' neural oscillations aligns consistently to the rhythmic input as young as
172 two months of age (Ní Choisdealbha et al., 2023).

173 Ní Choisdealbha and colleagues further demonstrated group-level convergence towards
174 a common *preferred phase angle* for rhythmic AV speech during the first year of life,

175 potentially reflecting convergence to an “optimal” phase angle for speech processing as
176 envisaged by Schroeder et al. (2008). In Ní Choisdealbha et al. (2023), the term “phase angle”
177 or “mean phase angle” referred to the angle of a cortical response at a given frequency at the
178 time at which it was observed, averaged over multiple observations. The same terminology is
179 used here. As the visual stimulus depicted an actor silently repeating a syllable every 500ms
180 (i.e. at a 2Hz rate), observations from the EEG time-frequency data were taken every 500ms.
181 Observations were taken across multiple EEG frequency bins for the purpose of statistical
182 comparison, but the sampling rate was always 500ms. "Preferred phase" thus refers to the
183 tendency for neural oscillations to occur around a particular angle (timepoint) relative to a
184 specific event at which it is measured (e.g. a repeated sound), and hence encompasses both the
185 mean phase angle of the observed oscillation and the consistency with which it occurs at or
186 around the same angle each time it is measured. Preferred phase can either be measured at an
187 individual level, by looking at an individual's repeated responses to a stimulus, or at a group
188 level, by examining whether the preferred phases of the individuals in the group fall into a
189 similar alignment. Group preferred phase may thus indicate the optimal alignment between a
190 stimulus and a neural response regarding the tracking of speech rhythm. Consequently, infants
191 whose phase angles are closer to this “preferred” alignment may have better language
192 outcomes. In our prior investigations with auditory-only and AV rhythmic speech, individual
193 differences in infant preferred phase indeed affected subsequent language development. Infant
194 preferred phase to auditory-only and AV speech was related to a range of measures of language
195 acquisition administered from 12 – 24 months (Ní Choisdealbha et al., 2023).

196 In the current report, we use a longitudinal approach to explore whether visual rhythmic
197 speech is processed by the infant brain. As noted, we adapted a rhythmic syllable repetition
198 EEG paradigm based on a 2Hz repetition rate, used previously with children in tests of TS
199 theory (Power et al., 2012, 2013; Keshavarzi et al. 2022a). This repetition rate corresponds to

200 the rate of the production of stressed syllables across languages (Dauer, 1983), and is within
201 the delta band of EEG. This is important in the context of TS theory, given that the phase
202 alignment of delta band rhythms to syllables differs in dyslexia, as does neural encoding of
203 delta-band speech envelope information (Power et al., 2012, 2013, 2016). The stressed syllable
204 rate has also been shown to be strongly represented in infant EEG responses to IDS (Attaheri
205 et al., 2022a; Menn et al., 2023). The Cambridge UK BabyRhythm study used one low-
206 frequency repetition rate (2Hz) in multiple modalities (speech, non-speech, auditory-only,
207 visual-only, AV). This was designed to enable estimation of cortical tracking of auditory-only,
208 visual-only and AV speech (Power et al., 2012; Ní Choisdealbha et al., 2022, 2023). In Power
209 et al.'s original (2012) study, typically-developing 13-year-old children showed significant
210 tracking of visual-only speech, and the phase coherence of cortical tracking across trials was
211 related to performance on a language task. These results contrast with those of Tan and
212 colleagues (2022), who did not find significant visual-only speech tracking for IDS at any age
213 in their study, which included four-year-olds and adults in addition to five-month-old infants.

214 Thus, simple rhythmic speech may provide an easier stimulus to align to than IDS, and
215 consequently may enable us to examine whether individual differences in phase alignment
216 might affect language outcomes. Further, Power et al. (2012) also reported that visually-
217 presented speech altered the preferred phase of *auditory* entrainment. This was examined by
218 subtracting visual-only EEG from AV responses (AV-V) in the syllable repetition paradigm
219 and comparing the resulting values to responses to auditory-only speech. The comparison
220 between AV-V and auditory-only speech conditions revealed a change in preferred phase at 4
221 Hz. Accordingly, comparing infant visual-only and AV responses in the rhythmic speech
222 paradigm may reveal information about changes in preferred phase in infants, changes that
223 could be relevant to the visual system preparing the auditory oscillators to be at the “optimal”
224 phase angle for speech processing (Schroeder, 2008).

225 The core research question explored here was whether individual differences in phase-
226 alignment to rhythmic visual-only speech exert an influence on later language development. A
227 number of research questions were of interest contributing to this primary question. First, on
228 the basis of our prior infant and child data, and the adult literature regarding oscillatory
229 mechanisms in speech processing, we expected to find significant cortical tracking of visual-
230 only rhythmic speech by infants. To investigate cortical tracking, we examined whether cortical
231 power at the stimulus frequency (2Hz) increased in response to the stimuli. We also measured
232 inter-trial coherence, which tells us whether the temporal alignment of neural oscillations at
233 2Hz was similar, relative to stimulus onset, across trials. Second, we were interested in the
234 presence of preferred phase both within and across infants. For group level preferred phase, we
235 investigated whether individuals' mean angle of alignment to the rhythmic stimuli coalesce
236 around a similar point at a group level. For individual level preferred phase, we investigated
237 whether individuals' mean phase angles differed when watching visual speech compared to a
238 resting condition with no visual stimulation. These measures enabled us to investigate potential
239 developmental changes in visual-only speech tracking and phase alignment between five and
240 eight months of age. Third, following Power et al. (2012), we investigated whether the
241 availability of visual speech information affects auditory cortical tracking. We did this by
242 subtracting visual-only responses from AV responses and examining stimulus-driven cortical
243 power and phase alignment. Following Power et al.'s results, we might expect to find a
244 common preferred phase across infants at 4 Hz. Finally, we were interested in whether
245 individual differences in infant phase angles (a) to visual-only speech and (b) between AV and
246 visual-only speech would predict infant-led and parent-reported measures of language
247 development taken at 12, 18 and 24 months of age. If the phase alignment mechanisms for
248 processing visual speech information are intrinsic to speech processing and support language
249 development (rather than being dependent on language comprehension) we would expect to

250 see cortical tracking of visual-only rhythmic speech in pre-verbal infants, as well as phase
251 alignment of visual-only neural responses to the stimuli. We would also expect to see relations
252 between phase alignment and later language performance.

253

254 **2. Methods**

255 ***2.1 Participants***

256 Participants were infants in the longitudinal Cambridge UK BabyRhythm Project. In total, 122
257 infants (65 male, 57 female) were recruited to the project and 113 participated through the
258 entire EEG data collection period, visiting the lab for assessments with the rhythmic speech
259 paradigm at two, five, six, eight and nine months. The five- and eight-month appointments
260 measured infant tracking of visual-only speech, at six and nine months an AV paradigm was
261 employed. For the visual speech assessments, 36 infants met the inclusion criteria of having at
262 least 5 trials in which they were detected by the eye-tracker as looking to the screen for at least
263 75% of the trial (19 infants present for both age groups). A strict looking time criterion was
264 necessary because when a child looks away from a visual-only stimulus, this interrupts
265 entrainment or other forms of oscillatory tracking because they are no longer receiving input
266 to track. For the AV-V assessments, there were 37 infants (six months) and 40 infants (nine
267 months) who watched at least five AV syllable trials for at least 75% of the trial data. However,
268 for these AV-V analyses, only infants who also had valid data for the preceding month's visual-
269 only paradigm appointment could be included. This left 18 infants with data at five and six
270 months, and 22 infants with data at eight and nine months. Full reasons for exclusions are given
271 in Appendix A, Table A1. Sample sizes for each analysis are given in Table A2 and Table A3.
272 Data collection occurred in the United Kingdom where data on race/ethnicity are protected and
273 cannot be collected without special cause.

274 An unknown technical factor affected mean phase angles that were computed for
275 datasets collected at appointments before or after February 14th 2018 (see Ní Choisdealbha et
276 al., 2023, for details). Due to the eye-tracking criteria, few infants tested before 14th February
277 2018 had valid data at five or eight months, (a minimum of two at eight months to a maximum
278 of eight at five months), so these infants were excluded. The final number of infants included
279 in the visual-only phase angle analysis was therefore 28 (five months), and 33 (eight months).
280 For the AV-V analyses, there were 15 infants with data at five and six months, and 21 at eight
281 and nine months.

282 ***2.2 Stimuli and Apparatus***

283 Stimuli were a video of a female speaker, recorded from the shoulders up, looking at the camera
284 and repeating the syllable “Ta” twice per second for 12 seconds. Videos were presented on a
285 TX300 eye-tracker screen (Tobii AB, Stockholm, Sweden) which recorded eye-tracking data
286 at a rate of 300 samples per second. At five and eight months, the stimuli were silent; at six and
287 nine months, the repeated syllable was played through 2020i speakers (Q Acoustics Ltd.,
288 Bishop’s Stortford, UK) driven by a Topaz AM5 stereo amplifier (Cambridge Audio Ltd.,
289 London, UK). The experiment was run using PsychToolbox (Brainard, 1997; Pelli, 1997;
290 Kleiner, Brainard & Pelli, 2007) via Matlab r2018b.

291 EEG data were recorded using an infant-friendly 64-channel Geodesic Sensor Net
292 (Electrical Geodesics, Inc., Eugene, OR, USA) connected to a GES 300 amplifier recording
293 via Netstation acquisition software at a sampling rate of 1kHz. A Cedrus StimTracker was used
294 to insert events into the EEG data to mark the onset of each “beat” of the syllable.

295 ***2.3 Procedure***

296 Recording took place in an EEG cubicle with infants sitting in an infant chair facing the eye-
297 tracker screen. With the EEG cap plugged in and recording, infants’ eye positions were
298 detected using the Tobii infant-friendly five-point calibration procedure. At five and eight

299 months, infants watched the stimulus videos consecutively. There was no defined inter-
300 stimulus interval but a brief buffering period of variable duration between each video, in which
301 no face was visible, meant that stimulation was not continuous. If the infant was no longer
302 engaged by the experimental stimuli on the screen, one or more attention-grabber stimuli of
303 varying lengths were played to re-engage their attention via novelty. At six and nine months
304 the syllable stimuli were also interspersed with other empirical stimuli showing a ball bouncing
305 on a drum and creating a 2Hz drumbeat sound (Ní Choisdealbha et al., 2022, 2023). Overall,
306 infants were shown 45.5 ± 11.9 (mean \pm SD) syllable videos at five months, 42.1 ± 9.7 at six
307 months, 50.7 ± 9.9 at eight months and 45.4 ± 6.3 at nine months. Once infants had seen
308 approximately 50 videos of the syllable or had become fussy or inattentive to the stimuli, EEG
309 recording took place in silence with no rhythmic visual or auditory stimulation for three to five
310 minutes, referred to hereafter as the “resting state”. During this part of the recording, a
311 researcher sat beside the infant and silently blew bubbles or pointed to images in a picture book
312 in a non-rhythmic way. This resting state condition was consistent across all recordings in the
313 Cambridge UK BabyRhythm project (see also Attaheri et al., 2022a; Ní Choisdealbha et al.,
314 2023).

315 The decision to present the AV and visual-only stimuli at different appointments was
316 made to limit attrition due to infant fussiness or tiredness, and maximise the quantity and
317 quality of data collected for each condition. While there may be age-related differences in
318 tracking and phase alignment, no age-related differences were found regarding infants'
319 responses to AV stimuli between 6 and 9 months (Ní Choisdealbha et al., 2023). Furthermore,
320 in the same sample, no age-related differences in stimulus tracking using the mTRF approach
321 were found between 7 and 11 months, though there were differences between 4 and 11 months
322 (Attaheri et al., 2022a,b). Given that our prior data indicates that cortical tracking and phase
323 alignment is similar between 6 and 9 months, we may assume that the month-gap necessary

324 here in recording infant responses to visual-only and AV speech did not exert strong effects on
325 the subtraction analyses.

326 ***2.4 EEG processing***

327 EEG data were preprocessed and analysed as described in Ní Choisdealbha and colleagues
328 (2022, 2023), using EEGLab 14 (Delorme & Makeig, 2004). Filtering, noise reduction (ASR;
329 Kothe & Makeig, 2013), bad channel detection and interpolation, epoching, baseline correction
330 and re-referencing were all performed as described in the previous manuscripts. Key details
331 are repeated in Appendix A.

332 Eye-tracking data were analysed using custom scripts, to determine for what proportion
333 of the trial the infant was looking at the screen. To account for brief losses in tracking, a
334 conservative interpolation procedure was implemented to fill in gaps in tracking of less than
335 50ms during which gaze did not move more than 10% of the width or height of the screen. If
336 the infant watched the screen for at least 75% of the length of the trial, that trial was included
337 in further analyses (see Table 1 for mean looking times, numbers of valid trials). The resting
338 state data was processed in the same manner as the syllable condition data, and epoched into
339 segments of the same length (see Table 1).

340 ***2.5 EEG analysis***

341 ***2.5.1 General approach***

342 The analysis approach mirrored that of other studies from the Cambridge UK BabyRhythm
343 project (Ní Choisdealbha et al., 2022; 2023). Starting with trial-level data, we (1) use a
344 normalised FFT to check if there is an increase in neural power in response to the stimulus at
345 the stimulus frequency, to establish whether tracking is taking place. We then (2) use an inter-
346 trial coherence (ITC) measure to see if changes in power are time-locked to the stimulus. This
347 tells us if the infant brain tracks the visual input precisely, or if the response is variable.
348 Regarding the individual “beats” (syllables or similarly-spaced instances within the resting

349 state segments), we then use the time-frequency data from step two to examine (3) the mean
350 angle of the 2Hz response recorded at the same time as each beat of the stimulus was played
351 (phase angle). This time-frequency data is also used to explore (4), the consistency at which
352 each "beat" of the response occurred at that angle (indexed by vector length). The angle of the
353 2Hz response to the visual speech stimulus in step (3) indicates whether the timing of that
354 response lines up differently in response to the stimulus, versus during resting state. The vector
355 length measure in step (4), like ITC, indexes consistency, but across beats rather than trials.

356 The final step (5) of EEG analysis moves from individual-level data to group-level data.
357 The individual analyses in steps three and four tell us whether, within infants, there are
358 differences in angle and consistency across conditions. The group-level preferred phase
359 analysis (step 5) tells us if there is a common preferred phase alignment between stimulus and
360 response, which in turn might suggest an "optimal" alignment – something which is important
361 for the language analyses. For each of these different EEG measures (relative power [1], ITC
362 [2], individual phase angle [3], individual vector length [4], and group phase angle [5]) we
363 examine both responses to the visual-only stimuli, and the values obtained when the response
364 to the visual-only stimulus is subtracted from the response to the AV stimulus.

365 Like other neural analyses within the Cambridge UK BabyRhythm project, we take a
366 whole-head approach to relative power and to ITC (Attaheri et al., 2022; Ní Choisdealbha et
367 al., 2022). The assumption is that, if tracking is taking place, the phenomenon should be robust
368 enough to pick up even when averaging across the scalp. For the results related to phase angle
369 ([3], [4] and [5] in the above list), it was possible that differences in phase profile across
370 different scalp regions of interest (ROIs) would attenuate locally-relevant results (for
371 depictions of phase angle differing by location, see Doelling et al., 2019; Ní Choisdealbha et
372 al., 2023). Consequently, we followed previous work (Cabral-Calderin & Henry, 2022; Ní
373 Choisdealbha et al., 2023) in using a combination of electrode-based testing and visual

374 inspection to identify electrodes with a stronger 2Hz response for the visual-only condition,
375 and a stronger 2Hz response for the AV – V data, and computing phase angle over the relevant
376 clusters. These included a mid-right frontal (MRF) and parietal (P) cluster for the visual-only
377 analyses, and right temporo-parietal (RTP), mid-right frontal (MRF), mid-left frontal (MLF),
378 and left temporal (LT) clusters for the AV – V analyses. Details including specific electrodes
379 can be found in Appendix B in the supplementary materials. Topographic plots showing where
380 the scalp-level 2Hz responses differ between visual stimulation and resting state, and between
381 audiovisual and visual stimulation, can be found in Figure B1 in the supplementary materials.

382 2.5.2 Relative power

383 For step (1) above, data were analysed using the same approaches as used by Ní Choisdealbha
384 et al. (2022; 2023). The relative power analysis indicates whether an increase in 2Hz power
385 took place. A Fast Fourier Transform (FFT) was run on the data using the MTMFFT approach
386 in Fieldtrip (Oostenveld et al., 2014) with a Hanning taper, using the 10.5-second length of the
387 trial and with frequency bins in 0.1Hz steps between 0.5 and 15Hz. Typically, power is greater
388 in the lower frequencies of the neural signal. To correct for this, data were then normalised to
389 neighbouring frequency bins by subtracting the mean power of the four nearest frequency bins
390 (two either side) from each frequency bin (Nozaradan et al., 2011). This means that results
391 show which frequency bins had greater (or lesser) power relative to its neighbouring bins.
392 Datapoints more than 3.5 times the standard deviation away from the mean were excluded as
393 outliers.

394 To obtain the difference values for the AV-V analyses, the normalised FFT values at
395 each frequency bin and in each condition (syllable and resting state) for the visual-only
396 paradigm were subtracted from the corresponding values for the AV paradigm at the next age-
397 point (i.e. five month values subtracted from six months, eight month values from nine
398 months).

399

400 2.5.3 Inter-trial coherence (ITC)

401 For step (2) above, to find whether changes in power were time-locked to the stimulus, a time-
402 frequency analysis was run in Fieldtrip, using the wavelet method with five Morlet wavelets
403 and taking 0.1Hz steps between 1 and 15Hz. The ITC analyses show whether the timing (or
404 phase) of the 2Hz response was consistent across trials. This would indicate that the increase
405 in 2Hz power, if present, was tagged to a particular event (i.e. the onset of the stimulus). The
406 ITC approach reflects both overall increases in power as well as phase consistency.
407 Consequently, the effects of increases in power, regardless of phase, were corrected for by
408 subtracting jittered from non-jittered ITC data so that the final ITC score would reflect *phase*
409 *consistency* across trials regardless of power (this analysis was also performed in Ní
410 Choisdealbha et al., 2022, and more information is available in Appendix A of the
411 Supplementary Materials). Datapoints more than 3.5 times the standard deviation away from
412 the mean were excluded as outliers. To obtain the difference values for the AV-V analyses, the
413 same subtraction approach as for the FFT was applied to the (non-jittered minus jittered) ITC
414 values.

415 2.5.4 Phase angle and vector length

416 The values for use in analysis steps (3), (4) and (5) were computed using the Circular Statistics
417 Toolbox (Berens, 2009). The approach used for the trial-level ITC data was repeated for each
418 of the ROIs individually, and the *circ_mean* function was used to compute the mean phase
419 angle and resultant vector lengths of infants' neural responses at each analysis frequency. Phase
420 angles were observed every 500ms from 1.5 to nine seconds after trial onset. Negative mean
421 phase angle values had two times pi added to them so that all results fell between 0 and 360
422 degrees.

423 For the phase angle and vector length differences, the mean phase data for each infant
424 at each age for each frequency bin for the visual-only speech was subtracted from the matched
425 responses for the AV speech (i.e. five month values subtracted from six months, and eight
426 months from nine months). The mean phase angle and vector length of this “phase difference”
427 for each infant was then obtained using the *circ_mean* function. In each case (visual-only or
428 AV-V) these values could then be compared by age, condition, ROI, and frequency bin to
429 examine their effects on within-individual preferred phase (steps 3 and 4) or compared to a
430 uniform distribution to examine group preferred phase (step 5; see sections 2.6.2. and 2.6.3
431 below).

432 **2.6 Statistical analysis**

433 2.6.1 Relative power and inter-trial coherence

434 Using lmerTest (Kuznetsova, Brockhoff & Christensen, 2017) in R, linear mixed effects
435 models were run on the normalised FFT and ITC values. There were fixed factors of age (five
436 versus eight months), condition (visual only versus resting state), frequency bin, and their
437 interactions. There was a random intercept on participant identity and where it was possible to
438 include a random slope on age without creating convergence or boundary issues, this is noted
439 in the results section. Satterthwaite-corrected ANOVAs were run on the models to determine
440 variables and interactions that contributed significantly to the models, while beta estimates are
441 reported to illustrate the size and direction of effects.

442 In all cases, the key interaction is that between frequency bin and condition, where an
443 increase in relative power and intertrial coherence at the 2Hz frequency bin in the syllable
444 condition would indicate that cortical tracking occurred at the stimulus frequency. The 2 Hz, 3
445 Hz, 4 Hz, 5 Hz, 6 Hz and 7 Hz frequency bins were included in the models, with 1 Hz also
446 included in the FFT model but not the ITC model (due the wavelet approach omitting values
447 on the edges of the time-frequency range). The inclusion of multiple frequency bins for

448 comparison purposes allows us to examine whether the increases in power and ITC occurred
449 at the 2Hz frequency bin specifically. While effects might be expected at the harmonic
450 frequencies (4Hz, 6Hz), an increase in relative power or ITC would not be expected at the other
451 frequencies. For comparability with related work (e.g. Attaheri et al., 2022a,b, Power et al.,
452 2013) we define 1Hz, 2Hz and 3Hz as delta, and 4Hz to 7Hz as theta.

453 For the relative power and ITC AV-V analyses, the same analyses were run on the AV-
454 V values. In this case the age factors were 5/6 months and 8/9 months; the conditions were the
455 AV-V difference and the difference in corresponding resting state values between the paired
456 recordings.

457 2.6.2 Vector length and phase angle

458 The analyses of individual-level phase angle and vector length used mixed effect regressions
459 with factors of age, condition, frequency bin and ROI, along with their interactions. Vector
460 length indexes whether an infant's neural responses at a given frequency occurred at the same
461 phase angle across trials, or if and by how much they varied. Analysis of individuals' mean
462 phase angle across conditions and frequency bins tell us whether the angle of phase alignment
463 to the stimulus differed from the mean phase angle observed at the sampled timepoints during
464 resting state. Phase angle is not a linear measure. High values are adjacent to low values (0° is
465 also 360°). To account for the circular nature of the phase angle data, models were run on the
466 sine and cosine values of each angle rather than the raw radian values. Sine and cosine provide
467 a linear representation of the circular features of the data (see also Ní Choisdealbha et al., 2023).

468 2.6.3 Group-level phase angle

469 The individual phase angle analyses tell us whether phase angle (or rather, its sine and cosine)
470 differed between conditions and frequency bins. The group-level analysis, conversely, tells us
471 whether responses converged across individuals, that is, whether infants showed a similar
472 preferred phase angle in response to the stimuli. The Rayleigh statistic was used to examine

473 whether mean phase angles on a *group* level converged toward a particular angle. It indicates
474 whether a set of circular values differ from a uniform distribution, that is, whether they are
475 clustered in a particular direction. The Rayleigh statistic for each group-level distribution of
476 phase angles were obtained using the *circ_rstat* function in the Circular Statistics Toolbox
477 (Berens, 2009). Based on prior work (Ní Choisdealbha et al., 2023), significant clustering of
478 responses was expected at 2Hz for the six- and nine-month responses to the AV syllable and
479 the same clustering was examined for the visual-only five- and eight-month responses.
480 Following Power and colleagues (2013), 4Hz clustering was examined for the AV-V analyses.

481 ***2.7 Language development analyses***

482 The language measures were decided *a priori* (Rocha et al., 2022 preprint, for rationale) and
483 were grouped into parent-reported standardised measures and infant-led experimental
484 measures encompassing phonology, communicative gesture, and vocabulary. The same
485 measures were used for all Cambridge UK BabyRhythm papers using neural predictors to
486 estimate language outcomes (Attaheri et al., 2022 preprint; Ní Choisdealbha et al., 2023). The
487 parent-reported measures were receptive and productive vocabulary scores in the Lincoln CDI
488 (Meints, Fletcher & Just, 2017) at 24 months. The infant-led measures were scores from three
489 tasks completed by the infants: (a) a binary measure of whether the infant could point
490 communicatively or not at 12 months; (b) their scores on the Computerised Comprehension
491 Task (CCT; Friend & Keplinger, 2003) at 18 months, in which infants point to pictures labelled
492 by the experimenter; and (c) their scores on a non-word repetition (NWR) task conducted at 24
493 months, encompassing the proportion of correctly-repeated consonants as well as accuracy in
494 reproducing the number of syllables and the stress pattern when repeating an item. Further
495 details on these tasks can be found in Rocha and colleagues (2022 preprint).

496 Relations between language outcome measures and the measures of phase alignment
497 were assessed using circular-linear correlations between phase angle and the outcome measures

498 (*circ_corrcl* in the Circular Statistics Toolbox, Berens, 2009), and linear-linear correlations
499 between vector length and the outcome measures. This follows the approaches of Power and
500 colleagues (2013) and Keshavarzi and colleagues (2022a) in their analyses of how preferred
501 phase relates to phonology in children with dyslexia. For all of the linear-linear correlations,
502 Pearson correlations were used, with the exception of the pointing measure, for which point
503 biserial correlations were run (Nagel, 2006). The alpha threshold used was adjusted based on
504 the number of measures in each group, that is, $\alpha = 0.01$ for the five infant-led measures, and α
505 = 0.025 for the two parent-estimated ones. To limit the number of tests performed, one ROI
506 was selected for analysis for each group of tests (visual-only and AV-V).

507 This approach diverges from our original plan for the visual-only speech data, which
508 was to replicate the analysis strategy used in our prior neural predictor papers (Attaheri et al.,
509 preprint; Ní Choisdealbha et al., 2023). In these papers, multivariate models were run using the
510 infant-led and parent-estimated measures, respectively, as dependent variables in each model.
511 The multivariate approach was also planned for the current data, with separate multivariate
512 models to be run on data from the five-month and eight-month age groups with predictor factors
513 of phase angle, vector length and the interaction between them. However, substantial
514 heteroscedasticity was observed when plotting the residuals of each model. Therefore, the
515 correlation-based approach described above was adopted. The AV-V sample sizes did not
516 provide sufficient degrees of freedom to run the multivariate models. Consequently, these data
517 were also analysed for each language measure individually using correlations.

518

519 **3. Results**

520 As there are many different analyses, we begin with a brief overview of the key results. First,
521 regarding *relative power* in response to *visual-only* speech, we do find greater relative power
522 at 2Hz in response to the visual-only stimulus, relative to the resting state and to all other

523 frequency bins. This indicates the presence of cortical tracking, and did not interact with age.
524 There were also differences in phase angle in response to the syllable relative to resting state
525 at 2Hz, both in terms of phase consistency (vector length) and the sine value of the angle
526 (predominantly over mid-right frontal electrodes).

527 Turning to the group-level phase analyses, the 2Hz mean phase angle does converge
528 for the 8-month-old infants in response to the visual-only syllable stimuli, but not for the 5-
529 month-old infants. The AV-V phase angles in response to the syllable at 8/9 months also
530 converge, at 4Hz. Accordingly, group preferred phases are present by 8 months.

531 For the other AV-V analyses, we find that power at 2Hz remains above the
532 comparison power value (taken from the resting state at the visual-only session subtracted
533 from that recorded at the AV session). The effect at 2Hz differs from the unrelated frequency
534 bins, but is not significantly different from that found at the harmonic bins (1Hz, 4Hz and
535 6Hz). There was no effect of age, and there were no effects on ITC. The crucial condition by
536 frequency interaction was not seen for the phase angle sine value, or vector length results,
537 and the cosine value showed only a trend ($p = 0.057$) towards a condition-by-frequency
538 interaction. The interaction arose because individuals' mean phase angles at 2Hz differed
539 between the syllable condition (AV-V value) and the resting state.

540 Finally, there were some significant relations between individual phase angle and
541 vector length measures, and subsequent scores on the language measures. The phase angle of
542 the response to the visual-only syllable at 8 months was related to parent-reported receptive
543 and productive vocabulary at 24 months.

544 ***3.1 Neural responses to visual-only speech***

545 3.1.1 Relative power

546 As noted in section 2.6.1, factors in the statistical analysis of relative power were condition
547 (visual only or resting state), frequency bin, age, and their interactions. There was a random

548 intercept on participant identity. We first examined the beat-corrected Fast Fourier Transform.
549 Our goal was to see if there was significantly greater relative power at 2Hz in response to the
550 visual stimuli, relative to the resting state EEG (Figure 1, a and b), and if this effect differed by
551 age. Table 2 provides the results of Satterthwaite-corrected tests of the fixed effects in each
552 model for the visual-only data (i.e. did each factor or interaction contribute significantly to each
553 model). Analysing effects of condition, frequency bin, age, and the interactions between them,
554 we found significant contributions of both frequency bin ($F(6, 867.88) = 4.12, p = 0.0004$) and
555 the frequency bin by condition interaction ($F(6, 867.88) = 3.32, p = 0.003$) to the model using
556 Satterthwaite-corrected ANOVAs. The estimates show that this was driven by a greater relative
557 power in the 2Hz frequency bin in response to the syllable, relative to the control condition
558 (2Hz relative power significantly higher than 7Hz base case, $\beta = 0.34, SE = 0.11, p = 0.002$;
559 2Hz relative power significantly higher than all other frequency bins 1 to 6Hz, all β s = 0.24 to
560 0.39, all p s < 0.05). Accordingly, significant cortical tracking was present. It did not differ by
561 age, as there was no significant contribution of age ($F(1, 360.3) = 0.015, p = 0.902$), and the
562 interactions between age and the other variables did not contribute to model fit (age by
563 condition, $F(1, 903.4) = 1.902, p = 0.168$; age by frequency bin, $F(6, 867.88) = 0.326, p =$
564 0.923 ; age by condition by frequency bin, $F(6, 867.88) = 0.514, p = 0.798$). As these model
565 contribution results suggest, there were no significant beta estimates relating to age nor to its
566 interactions. For this and all other analyses, full model results are available in Appendix B
567 (Table B1, Table B2, Table B3, and Table B4).

568 3.1.2 Inter-trial coherence

569 The inter-trial coherence analysis to explore whether the stimulus drew rhythmic cortical
570 responses into a consistent phase across trials revealed no significant effects of condition,
571 frequency bin or age (this model allowed for a random slope on age as well as the random
572 intercept on participant identity), nor of their interactions. Thus, although the FFT results

573 showed an increased cortical response to the stimuli, at the stimulus frequency (i.e. cortical
574 tracking), the ITC results suggest that the temporal consistency of this tracking varied across
575 trials.

576 3.1.3 Vector length and phase angle

577 The same factors – age, condition, frequency bin, and their interactions – were used in the
578 linear regressions on vector length and on phase angle (in the form of its sine and cosine
579 values), with the inclusion of ROI as well. Moving from trial-over-trial consistency to beat-by-
580 beat consistency, the vector lengths of infants’ mean phase angles in response to each “beat”
581 or instance of mouth opening did reveal a significant contribution of frequency bin ($F(5, 1524)$
582 $= 84.63, p < 0.0001$), which was reflected in longer 2Hz ($\beta = 0.069, SE = 0.031, p = 0.024$)
583 and shorter 3Hz ($\beta = -0.096, SE = 0.031, p = 0.002$) vectors relative to the 7Hz base case. There
584 was also a significant contribution of condition ($F(1, 1569) = 13.5, p = 0.0002$) and a significant
585 condition by frequency bin by ROI interaction ($F(1, 1569) = 2.565, p = 0.025$). The interaction
586 arose from a significantly longer 2Hz vector in response to the syllable for the MRF base case
587 ($\beta = 0.1, SE = 0.042, p = 0.017$), with a non-significant drop in consistency for the 2Hz response
588 to the syllable over the parietal ROI ($\beta = -0.098, SE = 0.059, p = 0.099$). Accordingly, on an
589 individual-by-individual basis, condition (visual-only versus resting state) did affect the
590 consistency of the 2Hz neural response.

591 Satterthwaite-corrected ANOVAs on phase angle (sine and cosine values) revealed a
592 significant effect of condition on the sine values ($F(1, 1377) = 6.67, p = 0.01$). However, the
593 only significant beta estimate was a three-way interaction between condition, ROI, and
594 frequency bin ($\beta = 0.696, SE = 0.35, p = 0.047$) – suggesting that the size of the difference in
595 the angle of individuals’ 2Hz phase response to the syllable (versus during resting state) differs
596 by scalp region. The cosine model had significant contributions of the condition by ROI,
597 frequency bin by ROI, and condition by age by ROI interactions (Table 2). However, the

598 critical frequency bin by condition interactions were all non-significant, as were the beta
599 estimates. Nonetheless, the sine effect shows that observed phase angle at 2Hz differed between
600 the resting state and syllable conditions.

601 *3.2 Group-level phase angle*

602 Group-level Rayleigh analyses of clustering of the phase angles suggested that the visual-only
603 speech drew responses into a similar phase angle across infants at 8 months. For resting state
604 data, we would not expect any clustering of the 2Hz phase angles given the absence of a
605 stimulus, and indeed no significant clustering was found (5 months, $n = 28$: MRF, $Z = 1.663$,
606 $p = 0.191$, P, $Z = 1.928$, $p = 0.146$; 8 months, $n = 33$: MRF, $Z = 0.333$, $p = 0.72$, P, $Z = 0.302$,
607 $p = 0.743$). By contrast, there was significant group-level clustering for the syllable at eight
608 months in the parietal ROI ($Z = 3.206$, $p = 0.039$, $n = 33$; Figure 2), although not in the MRF
609 ROI ($Z = 0.201$, $p = 0.82$) nor at five months (MRF: $Z = 0.0578$, $p = 0.565$, P: $Z = 1.565$, $p =$
610 0.21 , $n = 28$). These group-level results suggest that, by 8 months, viewing visual-only speech
611 is moving oscillatory responses towards a common preferred phase angle over parietal
612 electrodes. Although the non-significant Rayleigh tests suggest that individuals' preferred
613 phase differs over the MRF ROI, the vector length result for the syllable at 2Hz over this ROI
614 indicates that there is nonetheless a preferred phase within individuals.

615 For the AV-V data, the 4Hz band was analysed in accordance with the findings of
616 Power and colleagues (2013). The assumption is that the AV-V values represent the response
617 to an AV stimulus corrected for the response to the visual stimulus; thus the AV-V values
618 reflects the response to auditory stimulation alongside any effects related to audio-visual
619 integration. The only ROI showing group-level clustering was the MRF ROI, for the 8/9-month
620 values in the 4Hz syllable condition ($Z = 4.207$, $p = 0.013$; other 2Hz and 4Hz values reported
621 in Appendix B, Table B5). This suggests that for the 8/9-month-olds, auditory responses in the

622 presence of visual stimulation are drawn towards a common 4 Hz phase. This indicates that
 623 there is *auditory* phase alignment in the theta band while watching an audiovisual stimulus.

624 **3.3 Difference between AV and visual-only speech**

625 3.3.1 Relative power

626 For the AV-V FFT model (allowing for a random slope on age group), Satterthwaite-corrected
 627 ANOVAs on the model revealed significant factors of condition ($F(1, 467.68) = 3.9, p = 0.049$)
 628 and a condition by frequency bin interaction ($F(5, 668.67) = 2.28, p = 0.045$; see Table 3 for
 629 all factors and all models). Estimates reveal a significant increase in power in response to the
 630 syllable (relative to resting state) for the 2Hz base case ($\beta = 0.4, SE = 0.15, p = 0.009$). The
 631 2Hz syllable power was greater than at 3, 5 and 7 Hz (β s = -0.65, -0.44, -0.527, p 's = 0.003,
 632 0.043, and 0.016 respectively) but not than 1Hz ($\beta = -0.26, p = 0.24$), 4Hz ($\beta = -0.2, p = 0.37$)
 633 or 6Hz ($\beta = -0.36, p = 0.095$; Figure 1, c and d). There were no effects of age group. Figure 3
 634 depicts the 2Hz relative power values for each of the five-, six-, eight- and nine-month sessions,
 635 showing the difference between AV and visual responses to the speech stimulus. All simple
 636 effects are reported in Appendix B, Table B4. Although the 2Hz peak in Figure 1(c),
 637 corresponding the the 5/6-month AV-V, appears larger than that in Figure 1(d) (8/9-month AV-
 638 V), Figure 3 highlights why this age-related difference is non-significant in the model – a
 639 couple of infants showed very high relative 2Hz power at the younger ages (see Figure 3
 640 caption for note on outliers).

641 The ITC (allowing for random slope on age group), phase angle (sine and cosine) and
 642 vector length models were then run. The vector length, sine and cosine models all had
 643 significantly contributing factors according to the Satterthwaite-corrected ANOVA (see Table
 644 3). These significant factors were the frequency bin for vector length ($F(5, 1525) = 4.141, p =$
 645 0.001), the condition by age interaction for sine values ($F(1, 1545) = 6.512, p = 0.011$), and a
 646 marginal effect of condition for cosine values ($F(1, 1543) = 3.779, p = 0.052$ (see Table 3 for

647 all values). The crucial condition by frequency interaction did not make a significant
648 contribution to any model. Simple effects showed only a marginal interaction between the 2Hz
649 frequency and the syllable condition for the left temporal base case in the cosine model ($\beta =$
650 0.688 , $SE = 0.361$ $p = 0.057$).

651 *3.4 Language development analyses*

652 The results of the circular-linear and linear (Pearson or biserial) correlations with the
653 visual-only phase angle data are reported in Table 4. The MRF ROI was used in these analyses
654 as this is where we saw greater phase consistency within individuals. We note that there is a
655 well-documented right hemispheric preference for encoding rhythm (Sammler et al., 2015)
656 especially at slower rates (Vanvooren et al., 2014). Significant correlations were found between
657 phase angle in response to the stimulus at 8 months, and both receptive ($\rho(26) = 0.536$, $p =$
658 0.018) and productive vocabulary ($\rho(26) = 0.547$, $p = 0.015$) at 24 months, suggesting that
659 preferred phase may play a specific role in vocabulary development. The correlations with
660 vector length were all non-significant, suggesting that it may be the angle of phase alignment,
661 and not its consistency, which plays a role in the relation between visual speech processing and
662 later language acquisition.

663 The AV-V results relating to ROI were uninformative, hence we used the LT ROI for
664 the analyses regarding subsequent language development given that this region was used for
665 language analyses in related work (Ní Choisdealbha et al., 2023). Circular-linear and linear
666 (Pearson and biserial) correlations were run on the 5/6-month-olds' data and the 8/9-month-
667 olds' data. In these analyses, there were no significant associations for the infant-led measures
668 nor the parent-reported CDI scores. The closest trends related to the 5/6-month-olds' AV-V
669 vector lengths, and the 24-month NWR syllable-matching ($p = 0.038$) and stress-matching (p
670 $= 0.031$) measures, also the 8/9-month phase angle and 18-month CCT scores ($\rho(16) = 0.61$, p

671 = 0.037). Results for all measures are presented in Table 5 and should be interpreted in light of
672 the small sample sizes.

673 **Discussion**

674 Here we explored whether ‘visual prosody’ (low-frequency visual temporal information in
675 speech) evokes cortical tracking in five- and eight-month-old infants, whether individual
676 differences in preferred phase to visual-only speech affect subsequent language development,
677 and whether group preferred phase angles change when visual-only speech EEG is subtracted
678 from AV speech EEG (AV-V). The rhythmic speech paradigm employed here revealed that
679 infants show significant cortical tracking of visual-only speech at both five and eight months,
680 and that individuals do show phase consistency (longer phase angle vectors) in response to the
681 visual syllable. Convergence towards a group *preferred phase angle* for visual-only speech was
682 not significant at five months, but a preferred phase angle for visual-only speech had emerged
683 by eight months. Analyses of the difference in neural response between visual-only and AV
684 speech at 5/6 and 8/9 months revealed convergence towards a group *preferred phase angle* at
685 4 Hz rather than 2 Hz during cortical tracking, again for the older infants only. Consistent with
686 the adult and animal literature, therefore, the pre-verbal infant brain shows automatic use of
687 both visual and auditory prosodic information during speech processing. Further, individual
688 differences in these automatic neural mechanisms have consequences for language acquisition.
689 Each of these conclusions is discussed in greater detail below.

690 Regarding cortical tracking, the rhythmic speech paradigm employed here found
691 significant cortical tracking of visual-only speech at both ages tested, shown by the normalised
692 FFT results. The neural data from our 5-month-old infants contrasts with the negative result
693 for cortical tracking of visual-only speech for 5-month-old infants reported by Tan et al. (2022).
694 This could be due to some differences between the studies, such as the stimuli employed
695 (repeated syllables versus IDS) or attention thresholds (75% versus 15%). If we had utilised a

696 lower threshold in the current work, this would affect the subtractions between AV and visual-
697 only tracking, as an infant who disengages from an AV trial nonetheless continues to receive
698 rhythmic auditory input, whereas as infant who disengages from a visual-only trial no longer
699 receives any input. A final difference with the work by Tan et al. (2022) is that they used the
700 mTRF approach to correlate the neural response with the envelope of the presented speech,
701 while our focus was on individual and group-level *preferred phase*.

702 Regarding the phase analyses, we did not find significant ITC at either age, suggesting
703 that while cortical tracking occurred in response to the visual-only speech, the 2Hz neural
704 response was not drawn into a consistent alignment relative to stimulus onset for individual
705 infants across all trials. However, when we look at phase responses to individual stimulus beats,
706 broken down by region of interest, we see within-participant consistency in the vector lengths
707 of the 2Hz response to the visual syllable. The analyses of the circular data showed that phase
708 consistency was evident at the group level by 8 months. The group-level consistency
709 manifested over posterior, parietal electrodes while the individual-level consistency was
710 stronger over anterior, mid-right frontal electrodes. Activity at the scalp does not necessarily
711 reflect activity over adjacent cortex. However, it may be that the posterior response reflects the
712 evoked visual response, which is a fundamental sensory-perceptual response, while the mid-
713 right frontal response reflects other aspects of rhythmic processing. Accordingly, there may be
714 individual differences across infants even if consistency occurs within individuals.

715 As shown by these analyses, individual infants' mean 2Hz phase angles in response to
716 visual-only speech clustered in a specific direction, suggesting that there is a *preferred phase*
717 at which the infant brain processes 2Hz visual input, both individually and across the group.
718 Viewing visual-only speech thus affected the *timing* of the neural response, possibly by re-
719 setting the auditory cortex to the "optimal" state for processing anticipated succeeding
720 vocalisations (Schroeder et al., 2008). Regarding the apparently later emergence of group-

721 level preferred phase effects (older infants only) compared to cortical tracking effects (five and
722 eight months) it is notable that behavioural studies of infant speech processing using looking
723 measures or eye tracking measures indicate a change in preference between four and eight
724 months of age from looking at the eyes of a speaker to looking at the mouth (Lewkowicz &
725 Hansen-Tift, 2012; Pons, Bosch & Lewkowicz, 2015). The increased focus on the mouth as
726 infants get older is thought to provide infants with enhanced audiovisual cues during language
727 learning (Lewkowicz & Hansen-Tift, 2012). Such cues may include the visual prosody of the
728 speech directed to infants, who can detect such prosody from facial and head movements
729 (Kitamura et al., 2014). Accordingly, one speculative explanation for the significant group-
730 level preferred phase to visual-only speech found here at eight months is that it could reflect,
731 at least in part, behavioural changes in looking preference.

732 Turning to the AV-V analyses, the FFT statistics revealed significant tracking of the
733 rhythmic auditory stimulation in the presence of visual input. However, the individual-level
734 phase results were non-significant, suggesting that the alignment of neural response to stimulus
735 was variable. Nonetheless, the Rayleigh statistics showed that at 8/9 months (though not 5/6
736 months), 4 Hz AV-V angles converged towards a common direction, indicating a group-level
737 preferred phase. This AV-V result partially replicates earlier findings with older children.
738 Power and colleagues (2012) found 4 Hz responses to the rhythmic speech stimulus for AV
739 and auditory-only speech, but not visual-only speech. It is notable that in the previous analysis
740 of the infant AV-only speech data (reported by Ní Choisdealbha et al., 2022; these analyses
741 included additional six- and nine-month-old infants who did not meet the looking time
742 threshold used for the current paper), there was a significant increase in power at 4 Hz in
743 response to the AV syllable. Thus, auditory stimulation at 2 Hz – in the presence or absence of
744 visual stimulation – has effects at other, related frequencies, which visual stimulation alone
745 does not. A simple although speculative explanation could be that the auditory component of

746 AV speech contains 4 Hz information that is not carried by the visual input, but that is
747 potentially modulated by the visual input. One interpretation could be that for both infants and
748 13-year-olds (Power et al., 2012), 4 Hz AV-V effects are indicative of the re-setting of the
749 auditory cortex to the “optimal” state for processing the succeeding vocalisations by visual
750 speech information (Schroeder et al., 2008). Alternatively, the significant theta band responses
751 could reflect the automatic coupling of the auditory system to the motor system when speech
752 is heard, an effect indicated by adult auditory-only rhythmic speech data collected by Assaneo
753 and Poeppel (2018). Both of these speculative explanations could be examined in future
754 research.

755 Regarding subsequent language acquisition, if cortical tracking of visual prosody plays
756 a key role in neural speech processing, then infant responses to visual-only speech should be
757 predictive of later language outcomes. Individual differences in phase alignment to visual-only
758 speech significantly affected subsequent language development for some measures. At eight
759 months, individuals’ preferred phase to visual-only speech predicted parental estimates of both
760 receptive and productive vocabulary development at 24 months as measured by the UK-CDI
761 (Meints, Fletcher & Just, 2017). In prior analyses with these same infants, phase angle in
762 response to the AV syllable at nine months also significantly affected 24-month UK-CDI scores
763 (Ní Choisdealbha et al., 2023). If the role of visual speech is to promote phase-resetting for
764 optimal auditory processing, then we would expect AV-V phase responses – reflecting auditory
765 processes in the presence of visual speech – to also predict later language acquisition. We did
766 not find any robust relations between AV-V phase and language acquisition in this study.
767 However, some trends were present which could be explored in a larger sample. Particularly
768 interesting in the light of Power et al.’s (2013) data were those related to the NWR task, which
769 is a speech production measure, and can also be considered a measure of phonological

770 processing skills when used with very young children (Moritz et al., 2013; Kalashnikova et al.,
771 2021).

772 Our study has some limitations. There are myriad interesting questions about how the
773 tracking and processing of visual speech develops, including identifying the cortical sources of
774 such tracking in auditory and visual areas of the brain, and whether purported phase alignment
775 mechanisms are specific to visual speech or apply to other forms of visual or AV stimuli. These
776 questions were beyond the scope of the current study. A second limitation is that, despite a
777 large longitudinal EEG sample of 113 infants, there was substantial participant attrition. This
778 was primarily driven by the requirement for participants to have clean, usable EEG data for at
779 least five trials in which they attended to the screen for at least 75% of a 12-second video.
780 Furthermore, for the AV-V analyses, infants had to meet this criteria for two consecutive
781 testing sessions. Such attrition is a common finding in infant studies (see discussion in Tan et
782 al., 2022). Nevertheless our EEG analysis sample sizes of 18 (5/6-month AV-V) to 33 (eight-
783 month visual only) are representative of the infant literature (e.g. there were 18 five-month-
784 olds in Tan et al.). Given the parallels between our findings regarding preferred phase and those
785 reported by Power et al. (2012), the results obtained with the rhythmic repetition paradigm may
786 be treated as robust despite reduced sample sizes. A third limitation is that sample size was
787 particularly affected in the analyses predicting language outcomes, for which EEG and
788 language data from multiple testing sessions per infant was required. Nonetheless, many of the
789 effects of phase angle and vector length on later language performance reported here are
790 comparable to effects from the AV and auditory-only paradigms used in other test sessions
791 with the same infants, which employed larger samples (>100 infants, Ní Choisdealbha et al.,
792 2023). A final limitation is that the rhythmic speech paradigm is not comparable to connected
793 speech, although it shares some key acoustic parameters with IDS (Leong et al., 2017). While
794 the rhythmic speech paradigm is well-suited to phase analyses, of central interest here, it would

795 also be interesting to examine phase alignment during visual-only and AV-V speech processing
796 by infants when IDS is the input. In summary, our results show that neural oscillatory
797 mechanisms revealed by adult language processing studies, such as visual-only cortical
798 tracking (Park et al., 2016), appear to be engaged during visual speech processing in the first
799 year of life. These neural mechanisms are also engaged during auditory processing in the
800 presence of visual speech, with cross-modal modulation effects in evidence, and predict
801 language outcomes.

802

803 *Conclusions*

804 Adult studies of speech processing have suggested that auditory and visual low-frequency
805 neuronal oscillations act in concert when human speech is the input, with cortical tracking of
806 both auditory and visual prosody enhancing speech comprehension. The rhythmic speech
807 paradigm used here with infants reveals some similar visual and auditory oscillatory
808 mechanisms in the first year of life, when the infant is still pre-verbal. Whether the
809 underlying neural mechanisms measured here are general sensory-perceptual mechanisms or
810 specific to speech and language is not a question that we address here. Instead, we show that
811 individual differences in the angle and consistency of phase responses to the rhythmic
812 syllable stimuli were associated with certain measures of language development.

813 Consequently, individual differences in the engagement of these neural mechanisms by pre-
814 verbal infants, which are automatic mechanisms and not under conscious control, appears to
815 have potential implications for the early processing and development of language.

816 Accordingly, some of the neural mechanisms revealed by adult studies may be intrinsic to the
817 neural speech processing architecture, rather than dependent on learning to comprehend and
818 produce speech.

819 **References**

- 820 Assaneo, M. F., & Poeppel, D. (2018). The coupling between auditory and motor cortices is
821 rate-restricted: Evidence for an intrinsic speech-motor rhythm. *Science*
822 *Advances*, 4(2), aao3842.
- 823 Attaheri, A., Ní Choisdealbha, Á., Di Liberto, G. M., Rocha, S., Brusini, P., Mead, N.,
824 Olawole-Scott, H., Boutris, P., Gibbon, S., Williams, I., & Goswami, U. (2022a).
825 Delta-and theta-band cortical tracking and phase-amplitude coupling to sung speech
826 by infants. *NeuroImage*, 247, 118698.
- 827 Attaheri, A., Panayiotou, D., Philips, A., Ní Choisdealbha, Á., Di Liberto, G.M., Rocha, S.,
828 Brusini, P., Mead, N., Flanagan, S., Olawole-Scott, H., & Goswami, U. (2022b).
829 Cortical tracking of sung speech in adults vs infants: A Developmental Analysis.
830 *Frontiers in Neuroscience*, 16, 842447.
- 831 Attaheri, A., Ní Choisdealbha, Á., Rocha, S., Brusini, P., Di Liberto, G. M., Mead, N., ... &
832 Goswami, U. (2022 preprint). Infant low-frequency EEG cortical power, cortical
833 tracking and phase-amplitude coupling predicts language a year later. *bioRxiv*, 2022-
834 11.
- 835 Berens, P. (2009). CircStat: a MATLAB toolbox for circular statistics. *Journal of Statistical*
836 *Software*, 31, 1-21.
- 837 Bergelson, E., & Swingley, D. (2012). At 6–9 months, human infants know the meanings of
838 many common nouns. *Proceedings of the National Academy of Sciences*, 109(9),
839 3253-3258.
- 840 Bourguignon, M., Baart, M., Kapnoula, E. C., & Molinaro, N. (2020). Lip-reading enables
841 the brain to synthesize auditory features of unknown silent speech. *Journal of*
842 *Neuroscience*, 40(5), 1053-1065.
- 843 Brainard, D. H. (1997). The Psychophysics Toolbox. *Spatial Vision*, 10, 433-436.

- 844 Cabral-Calderin, Y. & Henry, M.J. (2022). Reliability of neural entrainment in the human
845 auditory system. *Journal of Neuroscience*, 42 (5), 894-908.
- 846 Chandrasekaran, C., Trubanova, A., Stillittano, S., Caplier, A., & Ghazanfar, A.A. (2009).
847 The natural statistics of audiovisual speech. *PLoS Computational Biology*, 5(7),
848 e1000436.
- 849 Dauer, R. M. (1983). Stress-timing and syllable-timing reanalyzed. *Journal of Phonetics*,
850 11(1), 51–62.
- 851 Delorme, A., & Makeig, S. (2004). EEGLAB: an open source toolbox for analysis of single-
852 trial EEG dynamics including independent component analysis. *Journal of*
853 *Neuroscience Methods*, 134(1), 9-21.
- 854 Doelling, K. B., Assaneo, M. F., Bevilacqua, D., Pesaran, B., & Poeppel, D. (2019). An
855 oscillator model better predicts cortical entrainment to music. *Proceedings of the*
856 *National Academy of Sciences*, 116(20), 10113-10121.
- 857 Friend, M., & Keplinger, M. (2003). An infant-based assessment of early lexicon acquisition.
858 *Behavioral Research Methods*, 25, 302-309.
- 859 Giraud, A. L., & Poeppel, D. (2012). Cortical oscillations and speech processing: emerging
860 computational principles and operations. *Nature neuroscience*, 15(4), 511-517.
- 861 Goswami, U. (2019). Speech rhythm and language acquisition: an amplitude modulation
862 phase hierarchy perspective. *Annals of the New York Academy of Sciences*, 1453(1),
863 67-78.
- 864 Goswami, U. (2011). A temporal sampling framework for developmental dyslexia. *Trends in*
865 *Cognitive Sciences*, 15(1), 3-10.
- 866 Grant, K.W., & Braida, L.D. (1991). Evaluating the articulation index for auditory–visual
867 input. *Journal of the Acoustical Society of America*, 89, 2952-2960.

- 868 Grant K.W., & Seitz, P-F. (2000). The use of visible speech cues for improving auditory
869 detection of spoken sentences. *Journal of the Acoustical Society of America*, 108,
870 1197-1208.
- 871 Kalashnikova, M., Peter, V., Di Liberto, G. M., Lalor, E. C., & Burnham, D. (2018). Infant-
872 directed speech facilitates seven-month-old infants' cortical tracking of
873 speech. *Scientific Reports*, 8(1), 1-8.
- 874 Kalashnikova, M., Burnham, D., & Goswami, U. (2021). Rhythm discrimination and
875 metronome tapping in 4-year-old children at risk for developmental
876 dyslexia. *Cognitive Development*, 60, 101129.
- 877 Keshavarzi, M., Mandke, K., Macfarlane, A., Parvez, L., Gabrielczyk, F., Wilson, A., &
878 Goswami, U. (2022a). Atypical delta-band phase consistency and atypical preferred
879 phase in children with dyslexia during neural entrainment to rhythmic audio-visual
880 speech. *NeuroImage: Clinical*, 35, 103054.
- 881 Keshavarzi, M., Mandke, K., Macfarlane, A., Parvez, L., Gabrielczyk, F., Wilson, A.,
882 Flanagan, S., Goswami, U. (2022b). Decoding of speech information using EEG in
883 children with dyslexia: Less accurate low-frequency representations of speech, not
884 "Noisy" representations. *Brain and Language*, 235, 105198.
- 885 Kitamura, C., Guellai, B., & Kim, J. (2014). Motherese by eye and ear: Infants perceive
886 visual prosody in point-line displays of talking heads. *PLoS One*, 9(10), e111467.
- 887 Kleiner, M., Brainard, D., & Pelli, D. (2007). What's new in Psychtoolbox-3? *Perception*, 36,
888 ECVF Abstract Supplement.
- 889 Kothe, C. A., & Makeig, S. (2013). BCILAB: a platform for brain-computer interface
890 development. *Journal of Neural Engineering*, 10(5), 056014.
- 891 Kuhl, P. K., & Meltzoff, A. N. (1982). The bimodal perception of speech in infancy. *Science*
892 218, 1138-1141.

- 893 Kuznetsova, A., Brockhoff, P.B., & Christensen, R.H.B. (2017). lmerTest Package: Tests in
894 linear effects models. *Journal of Statistical Software*, 82(13), 1-26.
- 895 Leong, V., Kalashnikova, M., Burnham, D., & Goswami, U. (2017). The temporal
896 modulation structure of infant-directed speech. *Open Mind*, 1, 78-90.
- 897 Lewkowicz, D. J., & Hansen-Tift, A. M. (2012). Infants deploy selective attention to the
898 mouth of a talking face when learning speech. *Proceedings of the National Academy
899 of Sciences*, 109(5), 1431-1436.
- 900 Luo, H., Liu, Z., & Poeppel, D. (2010). Auditory cortex tracks both auditory and visual
901 stimulus dynamics using low-frequency neuronal phase modulation. *PLoS Biology*, 8,
902 e1000445.
- 903 Mehler, J., Jusczyk, P., Lambertz, G., Halsted, N., Bertoni, J., and Amiel-Tison, C. (1988).
904 A precursor of language acquisition in young infants. *Cognition*, 29, 143-178.
- 905 Meints, K., Fletcher, K., & Just, J. (2017). The Lincoln Toddler Communicative
906 Development Inventory - A UK adaptation of the MacArthur-Bates Communicative
907 Development Inventory: Words and sentences (Toddler Form). [https://cpb-eu-
908 w2.wpmucdn.com/blogs.lincoln.ac.uk/dist/b/6736/files/2017/11/Lincoln_toddler_cdiv
909 2-2.pdf](https://cpb-eu-w2.wpmucdn.com/blogs.lincoln.ac.uk/dist/b/6736/files/2017/11/Lincoln_toddler_cdiv-2-2.pdf). Accessed 18/11/2022.
- 910 Menn, K. H., Michel, C., Meyer, L., Hoehl, S., & Männel, C. (2022). Natural infant-directed
911 speech facilitates neural tracking of prosody. *NeuroImage*, 251, 118991.
- 912 Moritz, C., Yampolsky, S., Papadelis, G., Thomson, J., & Wolf, M. (2013). Links between
913 early rhythm skills, musical training and phonological awareness. *Reading & Writing*,
914 26, 739 – 769.
- 915 Munhall, K.G., Jones, J.A., Callan, D.E., Kuratate, T., & Vatikiotis-Bateson, E. (2004).
916 Visual prosody and speech intelligibility: head movement improves auditory speech
917 perception. *Psychological Science*, 15, 133-137.

- 918 Munhall, K.G., Kross, C., & Vatikiotis-Bateson, E. (2002). Audiovisual perception of band-
919 pass filtered faces. *Journal of the Acoustical Society of Japan*, 21, 519-520.
- 920 Nagel, F. (2006). Point biserial correlation. *MATLAB Central File Exchange*.
- 921 Ní Choisdealbha, Á., Attaheri, A., Rocha, S., Mead, N., Olawole-Scott, H., Brusini, P.,
922 Gibbon, S., Boutris, P., Hines, D., Grey, C., Flanagan, S., & Goswami, U. (2022).
923 Cortical oscillations in pre-verbal infants track rhythmic speech and non-speech
924 stimuli. In Y. Gong & F. Kpogo, (Eds.), *BUCLD 46: Proceedings of the 46th annual*
925 *Boston University Conference on Language Development: Volume 2* (pp. 574-585).
926 Cascadilla Press.
- 927 Ní Choisdealbha, Á., Attaheri, A., Rocha, S., Mead, N., Olawole-Scott, H., Brusini, P., ... &
928 Goswami, U. (2023). Neural phase angle from two months when tracking speech and
929 non-speech rhythm linked to language performance from 12 to 24 months. *Brain and*
930 *Language*, 243, 105301.
- 931 Nozaradan, S., Peretz, I., Missal, M., & Mouraux, A. (2011). Tagging the neuronal
932 entrainment to beat and meter. *Journal of Neuroscience*, 31(28), 10234-10240.
- 933 Oostenveld, R., Fries, P., Maris, E., Schoffelen, & J.-M. (2011). FieldTrip: Open source
934 software for advanced analysis of MEG, EEG, and invasive electrophysiological data.
935 *Computational Intelligence and Neuroscience*, 2011, 156869.
- 936 Ortiz Barajas, M. C., Guevara, R., & Gervain, J. (2021). The origins and development of
937 speech envelope tracking during the first months of life. *Developmental Cognitive*
938 *Neuroscience*, 48, 100915.
- 939 Park, H., Kayser, C., Thut, G., & Gross, J. (2016). Lip movements entrain the observers' low-
940 frequency brain oscillations to facilitate speech intelligibility. *eLife*, 5, e14521.
- 941 Pelli, D. G. (1997). The VideoToolbox software for visual psychophysics: Transforming
942 numbers into movies, *Spatial Vision*, 10, 437-442.

- 943 Poeppel, D. (2003). The analysis of speech in different temporal integration windows:
944 cerebral lateralization as ‘asymmetric sampling in time’. *Speech Communication*, 41,
945 245-255.
- 946 Pons, F., Bosch, L., & Lewkowicz, D. J. (2015). Bilingualism modulates infants’ selective
947 attention to the mouth of a talking face. *Psychological Science*, 26(4), 490-498.
- 948 Pons, F., Lewkowicz, D.J., Soto-Faraco S., & Sebastián-Gallés, N. (2009). Narrowing of
949 intersensory speech perception in infancy. *PNAS*, 106, 10598-10602.
- 950 Power, A. J., Colling, L. J., Mead, N., Barnes, L., & Goswami, U. (2016). Neural encoding of
951 the speech envelope by children with developmental dyslexia. *Brain and Language*,
952 160, 1-10.
- 953 Power, A. J., Mead, N., Barnes, L., & Goswami, U. (2012). Neural entrainment to
954 rhythmically presented auditory, visual, and audio-visual speech in children. *Frontiers*
955 *in Psychology*, 3, 216.
- 956 Power, A. J., Mead, N., Barnes, L., & Goswami, U. (2013). Neural entrainment to rhythmic
957 speech in children with developmental dyslexia. *Frontiers in Human Neuroscience*, 7,
958 777.
- 959 Rocha, S., Ní Choisdealbha, Á., Attaheri, A., Mead, N., Olawole-Scott, H., Grey, C.,
960 Williams, I., Gibbon, S., Boutris, P., Brusini, P. & Goswami, U.C. (preprint).
961 Language acquisition in the longitudinal BabyRhythm cohort. psyarxiv.com/28c35
- 962 Sammler, D., Grosbras, M. H., Anwander, A., Bestelmeyer, P. E., & Belin, P. (2015). Dorsal
963 and ventral pathways for prosody. *Current Biology*, 25(23), 3079-3085.
- 964 Schroeder, C. E., Lakatos, P., Kajikawa, Y., Partan, S., and Puce, A. (2008). Neuronal
965 oscillations and visual amplification of speech. *Trends in Cognitive Science*, 12, 106-
966 113.

- 967 Sumbly W. H, Pollack I. 1954. Visual contribution to speech intelligibility in noise. *Journal*
968 *of the Acoustical Society of America*, 26, 212-215.
- 969 Tan, S. J., Kalashnikova, M., Di Liberto, G. M., Crosse, M. J., & Burnham, D. (2022). Seeing
970 a talking face matters: The relationship between cortical tracking of continuous
971 auditory-visual speech and gaze behaviour in infants, children and
972 adults. *NeuroImage*, 256, 119217.
- 973 Vanvooren, S., Poelmans, H., Hofmann, M., Ghesquiere, P., & Wouters, J. (2014).
974 Hemispheric asymmetry in auditory processing of speech envelope modulations in
975 prereading children. *Journal of Neuroscience*, 34(4), 1523-1529.
- 976 Walden, B.E., Prosek, R.A., Montgomery, A.A., Scherr, C.K., & Jones, C.J. (1977). Effects
977 of training on the visual recognition of consonants. *Journal of Speech and Hearing*
978 *Research*, 20, 130-145.
- 979 Yehia, H., Kuratate, T., & Vatikiotis-Bateson, E. (2002). Linking facial animation, head
980 motion and speech acoustics. *Journal of Phonetics*, 30, 555-568.
- 981 Yehia, H., Rubin, P., & Vatikiotis-Bateson, E. (1998). Quantitative association of vocal-tract
982 and facial behaviour. *Speech Communication*, 26, 23-43.
- 983
- 984
- 985
- 986
- 987
- 988
- 989
- 990
- 991

992

993

994 **Table 1:** Infant looking behaviour as measured by eye tracking at the different ages.

	5 months	6 months	8 months	9 months
Average % looking to each trial	38.49% (SD = 21.18%)	40.06% (SD = 15%)	38.05% (SD = 17.71%)	39.66% (SD = 14.67%)
Average trials watched for 75%+ of trial length	4.77 (SD = 6.21)	5.01 (SD = 5.12)	5.34 (SD = 6.46)	4.61 (SD = 4.91)
Average trials watched for 75%+ of trial length (babies with 5+ valid trials only)	11.53 (SD = 5.53)	10 (SD = 4.31)	12.69 (SD = 5.51)	9.5 (SD = 4.16)
Average number of resting state segments per infant who met the inclusion criteria (5+ syllable trials with 75%+ looking)	14.32 (SD = 8.56)	15.06 (SD = 6.62)	20.06 (SD = 5.62)	16.34 (SD = 6.24)

995

996

997 **Table 2:** Satterthwaite-corrected ANOVAs on model effects, neural responses to visual-only

998 speech

	FFT	ITC	PA: cosine	PA: sine	VL
--	-----	-----	------------	----------	----

Condition	F(1, 903.41) = 2.738#	F(1, 761.37) = 0.171	F(1, 1377) = 1.967	F(1, 1377) = 6.673**	F(1, 1569) = 13.5***
Frequency bin	F(6, 867.88) = 4.12***	F(5, 734.52) = 0.332	F(5, 1377) = 0.185	F(5, 1377) = 0.296	F(5, 1525) = 84.632***
Age	F(1, 360.3) = 0.015	F(1, 54.53) = 0.031	F(1, 447) = 0.125	F(1, 478) = 0.117	F(1, 39) = 0.002
ROI			F(1, 1377) = 1.803	F(1, 1377) = 0.227	F(1, 1525) = 1.539
Condition * Frequency	F(6, 867.88) = 3.318**	F(5, 734.52) = 0.362	F(5, 1377) = 0.376	F(5, 1377) = 0.241	<i>F(5, 1525) = 1.996#</i>
Condition * Age	F(1, 903.41) = 1.902	F(1, 761.41) = 0.0013	F(1, 1377) = 0.388	F(1, 1377) = 0.366	F(1, 1569) = 0.263
Condition * ROI			F(1, 1377) = 7.425**	F(1, 1377) = 0.045	F(1, 1525) = 0.487
Frequency * Age	F(6, 867.88) = 0.326	F(5, 734.52) = 0.81	F(5, 1377) = 0.274	F(5, 1377) = 0.121	F(5, 1525) = 0.137
Frequency * ROI			F(5, 1377) = 2.523*	F(5, 1377) = 1.346	F(5, 1525) = 0.979
Age * ROI			<i>F(1, 1377) = 2.724#</i>	F(1, 1377) = 0.058	F(1, 1525) = 0.004
Condition * Frequency * Age	F(6, 867.88) = 0.514	F(5, 734.52) = 0.025	F(5, 1377) = 1.415	F(5, 671) = 0.76	F(5, 1525) = 1.671
Condition * Frequency * ROI			F(5, 1377) = 0.681	F(5, 1377) = 0.761	F(5, 1525) = 2.565*
Condition * Age * ROI			F(1, 1377) = 11.162***	<i>F(1, 1377) = 3.809#</i>	F(1, 1525) = 0.093
Frequency * Age * ROI			F(5, 1377) = 1.143	F(5, 1377) = 0.492	F(5, 1525) = 0.314
Condition * Frequency * Age * ROI			F(5, 1377) = 0.478	F(5, 1377) = 0.941	F(5, 1525) = 0.237

999 # p < 0.1 * p < 0.05 ** p < 0.01 *** p < 0.001

1000 PA = phase angle; VL = vector length

1001

1002 **Table 3:** Satterthwaite-corrected ANOVAs on model effects, neural responses to

1003 audiovisual-minus-visual (AV-V) speech

	FFT	ITC	PA: cosine	PA: sine	VL
Condition	F(1, 467.68) = 3.895*	F(1, 419.15) = 2.798#	$F(1, 1543) = 3.779\#$	F(1, 1546) = 0.526	F(1, 1544) = 0.7573
Frequency bin	F(6, 446.91) = 1.579	F(5, 410.11) = 1.889#	F(5, 1522) = 0.489	F(5, 1526) = 0.441	F(5, 1525) = 4.141**
Age	F(1, 20.75) = 0.316	F(1, 31.09) = 1.048	F(1, 15) = 0.032	F(1, 20) = 1.158	F(1, 18) = 0.039
ROI			F(3, 1522) = 0.949	F(3, 1526) = 1.027	F(3, 1525) = 1.243
Condition * Frequency	F(6, 446.91) = 2.198*	F(5, 410.11) = 1.71	F(5, 1522) = 1.574	F(5, 1526) = 1.181	F(5, 1525) = 0.23
Condition * Age	F(1, 467.6) = 0.514	F(1, 419.44) = 0.876	F(1, 1543) = 0.934	F(1, 1546) = 6.512*	F(1, 1544) = 0.906
Condition * ROI			F(3, 1522) = 0.776	F(3, 1526) = 1.477	F(3, 1525) = 1.603
Frequency * Age	F(6, 446.91) = 0.233	F(5, 410.09) = 0.749	F(5, 1522) = 0.418	F(5, 1526) = 0.173	F(5, 1525) = 0.156
Frequency * ROI			F(15, 1522) = 0.749	F(15, 1526) = 0.435	F(15, 1525) = 0.826
Age * ROI			F(3, 1522) = 0.921	F(3, 1526) = 0.692	F(3, 1525) = 0.493
Condition * Frequency * Age	F(6, 446.91) = 0.304	F(5, 410.09) = 0.373	F(5, 1522) = 0.764	F(5, 1526) = 0.996	F(5, 1525) = 1.021
Condition * Frequency * ROI			F(15, 1522) = 0.92	F(15, 1526) = 0.75	F(15, 1525) = 1.105
Condition * Age * ROI			F(3, 1522) = 0.276	F(3, 1526) = 1.795	F(3, 1525) = 0.183
Frequency * Age * ROI			F(15, 1522) = 0.526	F(5, 1526) = 0.91	F(15, 1525) = 1.423
Condition * Frequency * Age * ROI			F(15, 1522) = 0.992	F(5, 1526) = 0.474	F(15, 1525) = 0.909

1004 # $p < 0.1$ * $p < 0.05$ ** $p < 0.01$ *** $p < 0.001$

1005 PA = phase angle; VL = vector length

1006

1007 **Table 4:** All circular-linear (phase angle) and linear (vector length) correlations between visual

1008 only phase responses over mid-right frontal ROI, and language measures. Alpha is corrected

1009 for multiple comparisons and is equal to 0.01 for the infant-led and 0.025 for the parent-
 1010 estimated measures, accounting for five different infant-led measures and two different parent-
 1011 reported vocabulary measures.

	Circular-linear correlations with phase angle		Linear correlations with vector length	
	5 months	8 months	5 months	8 months
Infant-led measures				
Pointing	$\rho(26) = 0.101,$ $p = 0.866$	$\rho(31) = 0.038, p$ $= 0.976$	$r(26) = 0.273, p$ $= 0.152$	$r(31) = 0.057, p =$ 0.75
CCT	$\rho(20) = 0.038,$ $p = 0.984$	$\rho(26) = 0.436, p$ $= 0.07$	$r(20) = 0.491,$ $p = 0.02$	$r(26) = -0.177, p =$ 0.555
Consonants correct (NWR)	$\rho(25) = 0.109,$ $p = 0.852$	$\rho(29) = 0.276, p$ $= 0.308$	$r(25) = 0.119, p$ $= 0.554$	$r(29) = -0.12, p =$ 0.521
Syllables correct (NWR)	$\rho(25) = 0.134,$ $p = 0.785$	$\rho(29) = 0.205, p$ $= 0.522$	$r(25) = 0.195, p$ $= 0.329$	$r(29) = -0.021, p =$ 0.911
Correct stress (NWR)	$\rho(25) = 0.234,$ $p = 0.477$	$\rho(29) = 0.174, p$ $= 0.627$	$r(25) = 0.043, p$ $= 0.831$	$r(29) = -0.126, p =$ 0.5
Parent-estimated measures				
Receptive vocabulary (CDI)	$\rho(24) = 0.25, p$ $= 0.443$	$\rho(26) = 0.536, p$ $= 0.018$	$r(24) = 0.199, p$ $= 0.331$	$r(26) = -0.334, p =$ 0.083
Productive vocabulary (CDI)	$\rho(24) = 0.102,$ $p = 0.874$	$\rho(26) = 0.547, p$ $= 0.015$	$r(24) = 0.232, p$ $= 0.255$	$r(26) = -0.174, p =$ 0.376

1012 Bold font: $p < 0.01$ (infant-led) or $p < 0.025$ (parent-estimated). Italic: $p < 0.05$.

1013

1014 **Table 5:** All circular-linear (phase angle) and linear (vector length) correlations between AV-
 1015 V phase angle/vector length over left temporal ROI, and language measures. Alpha is corrected
 1016 for multiple comparisons and is equal to 0.01 for the infant-led and 0.025 for the parent-
 1017 estimated measures, accounting for five different infant-led measures and two different parent-
 1018 reported vocabulary measures.

	Circular-linear correlations with phase angle		Linear correlations with vector length	
	5/6 months	8/9 months	5/6 months	8/9 months
Infant-led measures				
Pointing	$\rho(13) = 0.188,$ $p = 0.766$	$\rho(19) = 0.424, p$ $= 0.152$	$r(13) = 0.308, p$ $= 0.247$	$r(16) = -0.006, p =$ 0.98
CCT	$\rho(11) = 0.406,$ $p = 0.343$	$\rho(16) = 0.606, p$ $= 0.037$	$r(11) = 0.294, p$ $= 0.33$	$r(14) = 0.332, p =$ 0.209
Consonants correct (NWR)	$\rho(12) = 0.597,$ $p = 0.08$	$\rho(18) = 0.125, p$ $= 0.856$	$r(12) = 0.486, p$ $= 0.078$	$r(15) = -0.073, p =$ 0.781
Syllables correct (NWR)	$\rho(12) = 0.63, p$ $= 0.062$	$\rho(18) = 0.046, p$ $= 0.979$	$r(12) = 0.557,$ $p = 0.038$	$r(15) = -0.045, p =$ 0.865
Correct stress (NWR)	$\rho(12) = 0.629,$ $p = 0.063$	$\rho(18) = 0.342, p$ $= 0.31$	$r(12) = 0.576,$ $p = 0.031$	$r(15) = -0.315, p =$ 0.218
Parent-estimated measures				
Receptive vocabulary (CDI)	$\rho(12) = 0.521,$ $p = 0.149$	$\rho(16) = 0.125, p$ $= 0.869$	$r(12) = 0.286, p$ $= 0.322$	$r(13) = 0.153, p =$ 0.587

Productive vocabulary (CDI)	$\rho(12) = 0.587,$ $p = 0.09$	$\rho(16) = 0.244, p$ $= 0.586$	$r(12) = 0.342, p$ $= 0.232$	$r(13) = 0.229, p =$ 0.411
-----------------------------------	-----------------------------------	------------------------------------	---------------------------------	---------------------------------

1019 Bold font: $p < 0.01$ (infant-led) or $p < 0.025$ (parent-estimated). Italic: $p < 0.05$.

1020

1021

1022

1023

1024

1025

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037 **Figure captions:**

1038

1039 **Figure 1:** Relative power values for visual speech, resting state and AV-V conditions, across

1040 frequency bins. Relative power during the syllable stimulus is shown in blue and during the

1041 resting state in red. Panel (a) shows group mean relative power (shaded area indicates

1042 standard error) for the five-month visual speech data; panel (b) for the eight-month visual
1043 speech data; panel (c) plots the difference in relative power between the AV speech at six
1044 months and the visual-only speech at five months; and panel (d) plots the difference in
1045 relative power between the AV speech at nine months and the visual-only speech at eight
1046 months. Analyses related to panels (a) and (b) are in section 3.1.1; those related to (c) and (d)
1047 are in 3.3.1.

1048

1049 **Figure 2:** Circular plots showing individual mean phase angles and vector length in blue and
1050 group means in yellow for the visual-only condition and the AV-V condition.

1051

1052 **Figure 3:** Relative power at 2Hz, plotted by age group and condition. The 5mo (five-month)
1053 and 8mo (8-month) data relates to the visual-only speech; the 6mo (six-month) and 9mo (nine-
1054 month) data is from the AV speech. Each dot is an individual infant's mean relative power at
1055 2Hz during the resting state (red, left) or syllable stimulus (blue, right) (Outliers were already
1056 excluded for AV and VO separately, with no 2Hz values reaching outlier exclusion criteria,
1057 consequently no further exclusions were performed).

1058

1059 **Figure 4:** Significant relationships between phase angle or vector length in visual speech or
1060 AV-V analyses, and later language measures. Panel A depicts 2Hz phase angle in response to
1061 visual-only speech at 8 months and receptive vocabulary on the CDI; panel B depicts 2Hz
1062 phase angle in response to visual-only speech at 8 months and productive CDI vocabulary. In
1063 panels A and B, the blue dots are participants, ordered from the participant with the smallest
1064 phase angle to the largest (expressed in radians). Values are shown on the left y-axis. The red
1065 dots are language outcome measures, with values shown on the right y-axis. The red line is
1066 plotted by the `geom_smooth` function in `ggplot2` using the loess method to illustrate the trend

1067 in the data points. The sinusoidal aspect of the trend lines is indicative of the underlying
1068 circular-linear relationship to the phase angle data.