



HAL
open science

Atypical Audio-Visual Neural Synchrony and Speech Processing in children with Autism Spectrum Disorder

Xiaoyue Wang, Sophie Bouton, Nada Kojovic, Anne-Lise Giraud, Marie Schaer

► **To cite this version:**

Xiaoyue Wang, Sophie Bouton, Nada Kojovic, Anne-Lise Giraud, Marie Schaer. Atypical Audio-Visual Neural Synchrony and Speech Processing in children with Autism Spectrum Disorder. 2024. pasteur-04627417v1

HAL Id: pasteur-04627417

<https://pasteur.hal.science/pasteur-04627417v1>

Preprint submitted on 27 Jun 2024 (v1), last revised 24 Feb 2025 (v2)

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons Attribution - NonCommercial - NoDerivatives 4.0 International License

1 **Atypical Audio-Visual Neural Synchrony and Speech Processing in** 2 **children with Autism Spectrum Disorder**

3
4 Xiaoyue WANG ^{1,2*}, Sophie BOUTON ², Nada KOJOVIC ³, Anne-Lise GIRAUD ^{1,2†}, Marie
5 SCHAER ^{3†}

6 *1 Auditory Language Group, Department of Basic Neuroscience, University of Geneva, Geneva,*
7 *Switzerland*

8 *2 Institut Pasteur, Université Paris Cité, Hearing Institute, Paris, France*

9 *3 Autism Brain & Behavior Lab, Department of Psychiatry, University of Geneva, Geneva,*
10 *Switzerland*

11 [†]*Co-last authors*

12 *Correspondence: xiaoyue.wang@pasteur.fr; wangxiaoyue43@gmail.com

13
14 Running title: Atypical audio-visual synchrony in ASD

17 **Acknowledgment**

18 **Author Contributions:** Data acquisition and clinical resources: M.S., N. K.; Study design: X.W.,
19 N. K., M.S. and A-L.G.; Data analysis: X.W., S.B.; Writing-original draft: X.W.; Writing-review-
20 editing: X.W., S.B., N. K., A-L.G. and M.S.; Supervision: A-L.G., and M.S.

21 **Funding:** This work is supported by grants from the Swiss National Science Foundation (#163859,
22 #190084, #202235 & #212653 to M.S), by the National Centres of Competence in Research
23 (NCCR) Synapsy (Grant No. 51NF40–185897 to M.S.) and Evolving Language (Grant No.
24 51NF40_180888 to A-L.G), by grant from Fondation pour l’Audition (FPA IDA11 to A-L.G) as
25 well as by support from the Fondation Privée des Hôpitaux Universitaires de Genève
26 (<https://www.fondationhug.org>), and the Fondation Pôle Autisme ([https://www.pole-
27 autisme.ch](https://www.pole-autisme.ch)).

28 **Conflict of interest:** The authors declare no conflict of interest.

29 **Data and materials availability:** The unprocessed datasets for this manuscript are not publicly
30 available yet, because analysis of these data is ongoing as part of a longitudinal study, and results
31 are expected to be published in the future. When all data has been published, requests to access
32 the datasets should be directed to Dr Marie Schaeer, marie.schaer@unige.ch, the custom MATLAB
33 analysis scripts will be made available upon request to the Lead Contact, Xiaoyue Wang
34 (xiaoyue.wang@pasteur.fr).

35

36

37

1 **Abstract (240 words)**

2

3 Background: Children with Autism Spectrum Disorders (ASD) often exhibit communication
4 difficulties that may stem from basic auditory temporal integration impairment but also be
5 aggravated by an audio-visual integration deficit, resulting in a lack of interest in face-to-face
6 communication. This study addresses whether speech processing anomalies in young (mean age
7 3.09-year-old) children with ASD are associated with alterations of audio-visual temporal
8 integration.

9 Methods: We used high-density electroencephalography (HD-EEG) and eye tracking to record
10 brain activity and gaze patterns in 31 children (6 females) with ASD and 33 typically developing
11 (TD) children (11 females), while they watched cartoon videos. Neural responses to temporal
12 audio-visual stimuli were analyzed using Temporal Response Functions model and phase analyses
13 for audiovisual temporal coordination.

14 Results: The reconstructability of speech signals from auditory responses was reduced in children
15 with ASD compared to controls, but despite more restricted gaze patterns in ASD it was similar
16 for visual responses in both groups. Speech reception was most strongly affected when visual
17 speech information was also present, an interference that was not seen in TD children. These
18 differences were associated with a broader phase angle distribution (exceeding $\pi/2$) in the EEG
19 theta range in autistic children, signaling reduced reliability of audio-visual temporal alignment.

20 Conclusion: These findings show that speech processing anomalies in ASD do not stand alone and
21 that they are associated already at a very early development stage with audio-visual imbalance
22 with lousier auditory response encoding and disrupted audio-visual temporal coordination.

23

24 **Keywords (6 keywords):** Autism, Gaze direction, Speech envelope, Visual motion,
25 Audio-visual, Oscillation Phase entrainment

26

27

28

1 Introduction

2 Newborns are immediately attracted to the human voice. In utero exposure to speech sounds
3 enables them to accurately discriminate speech sounds at birth (1–3). Since vision develops with
4 a delay relative to hearing, babies only progressively discover that vocal stimuli are related to
5 facial movements. Unlike typically developing (TD) children, children with Autism Spectrum
6 Disorders (ASD) do not show this primary interest in speech (4–7). Instead, they tend to engage
7 in slow and repetitive visual exploration of their environment, which has been suggested to lead
8 to atypical interests over time (8–13). Focusing on visual aspects of their surroundings allows
9 children with ASD to explore the world at their own pace, keeping them away from highly dynamic
10 stimuli such as speech and biological motion (14–16), which are often perceived as overwhelming
11 (17,18).

12
13 A basic auditory dysfunction in ASD might lead to speech-processing anomalies that in turn
14 cascade into a decreased interest in speech (19–24). Atypical speech processing becomes apparent
15 very early in development, and the fact that early neural anomalies, such as delta, theta, and gamma
16 oscillations, accurately predict the severity of future language deficits could suggest that they are
17 causal to later difficulties in language comprehension and production (25–28). The tendency of
18 children with ASD to prefer static or slow visual processing (29–32) possibly exacerbates speech
19 reception challenges by counteracting dynamic audio-visual interaction, a crucial process for
20 speech reception in ecological (e.g., noisy) environments (33–36). Accordingly, exceedingly long
21 integration time windows for audio and visual stimuli have been reported in children with ASD
22 (31,32), implying disturbed integration of audio and visual stimuli in ASD.

23
24 Two essential mechanisms participate in audio-visual integration. The first one is the relative
25 timing of auditory and visual stimuli: when falling within approximately 250 ms of each other,
26 they are often perceived as a single event, potentially influencing each other (e.g. the McGurk
27 effect (37,38)). The second mechanism is the resynchronization provoked by the stimulus in one
28 sensory modality affecting neural responses in the other one (39–44). Orofacial visual movements
29 typically precede the onset of speech, leading to a resynchronization that sharpens the auditory
30 speech response (39,45). And independent from the integration/fusion of the exact visual and
31 speech content, visual resynchronization enhances speech processing by boosting the tracking of
32 the speech's syllabic structure. While audio-visual temporal integration anomalies in ASD are well
33 documented (31,32,46–48), audio-visual dynamic synchronization anomalies remain hypothetical.

34
35 Auditory and visual sensory processing both operate rhythmically (49–53). Visual speech
36 information (lip movements) is characterized by a dominant 2-7 Hz rhythm (theta band (54)) and
37 these quasiperiodic visual cues influence speech perception by modulating auditory neuronal
38 oscillations within the same theta range at about 5 Hz (53–58), corresponding to the typical audio-
39 visual (AV) integration temporal time around 250ms (39,61–64). The reset of auditory neural
40 oscillations triggered by visual input (65) rhythmically enhances auditory processing (66), a

1 phenomenon that is already observable in typical children (67). Despite the documented presence
2 of auditory processing anomalies in ASD around 3-year-old (24), we still ignore whether they are
3 associated with dynamic audio-visual synchrony anomalies.

4
5 This study fills this gap by investigating with high-density EEG the dynamics of auditory and
6 visual processing in young children with and without ASD, aged 1.13 to 5.56 years old, under
7 naturalistic audio-visual conditions, i.e. when children are watching a popular cartoon adapted to
8 their age. The goal is to compare the quality of the neural encoding/decoding of dynamic auditory
9 and visual stimuli and audio-visual temporal coordination across groups.

10 11 **Methods**

12 **1 Participants**

13 Participants were selected from the Geneva Autism Cohort, a longitudinal study that aims at better
14 understanding the developmental trajectories in young children with ASD. This cohort's protocol
15 has been detailed in previous studies (22,62,63). In this study, we used clinical and behavioral
16 assessments, as well as the electroencephalogram (EEG) recorded simultaneously with eye-
17 tracking when children were watching popular cartoon videos.

18
19 The sample comprised 31 children diagnosed with ASD (6 females, mean age = 3.09 years, SD =
20 0.91, age range: 1.74 - 5.14) and 32 TD peers (11 females, mean age = 2.95 years, SD = 1.31, age
21 range: 1.31 - 5.56). Selection criteria for all participants included: age below 6 years, data collected
22 during the participant's initial visit (i.e. at autism diagnosis for the autistic group), clear and
23 accurate markers associated with movie onset, usable raw data for four different movies, and focus
24 on the screen throughout all recordings. The age difference between the two groups was not
25 significant (Kolmogorov-Smirnov $D = 0.28$, $p = 0.17$).

26
27 ASD clinical diagnoses were meticulously confirmed using standardized tools: either the Autism
28 Diagnosis Observation Schedule-Generic (ADOS-G) (70) or the Autism Diagnosis Observation
29 Schedule, Second Edition (ADOS-2) (71). Recruitment of participants occurred through
30 specialized clinical centers and community-wide announcements. For the TD group, exclusion
31 criteria included any suspicion of atypical psychomotor development, a history of neurological or
32 psychological disorder, or having a first-degree relative with an autism diagnosis.

33
34 Informed consent was obtained from the parents of all participants prior to inclusion in the study.
35 The research was conducted with the ethical standards set forth by the Ethics Committee of the
36 Faculty of Medicine at the University of Geneva Hospital and adhered to the principles outlined
37 in the Declaration of Helsinki.

38 39 **2 Stimuli and Procedure**

1 To explore cortical processing of audio-visual stimuli, we employed a passive and naturalistic task
2 suitable for young children. This task involved viewing an age-appropriate French cartoon
3 "TROTRO" (72–75) (example: <https://www.youtube.com/watch?v=jT9C9WCIQr8&t=81s>). The
4 selection of "TROTRO" was based on its cognitive accessibility and appeal to the target age group.
5 Importantly, TROTRO, the main character, speaks and interacts verbally with other characters,
6 which allows us to isolate the speech soundtrack and associated visual motion and to probe related
7 brain responses. Participants watched four Trotro episodes, each lasting approximately 2.5
8 minutes. The videos were presented in a consistent, predetermined order to all participants. To
9 monitor the participant's visual engagement with the stimulus, Tobii Studio (Tobii® Technology,
10 Sweden) was used. The screen for the video display was configured with dimensions of 1200
11 pixels in height (29°38' visual angle) and 1920 pixels in width (45°53'), with a refresh rate of 60
12 Hz. This setup was optimized for clear and comfortable viewing in children. Participants were
13 seated at an optimal distance of approximately 60 cm from the screen (Figure 1A).

14

15 **2.1 Eye-tracking acquisition and analysis**

16 Gaze data were collected using the Tobii TX300 eye-tracking system (<https://www.tobii.com>),
17 which operates at a sampling rate of 300 Hz. This high-frequency data collection was instrumental
18 in assessing participants' visual exploration patterns during the cartoon viewing. The cartoon was
19 displayed in a frame that provided a visual angle of 26°47'(height) × 45°53'(width). Calibration
20 was performed using a child-friendly procedure integrated into the Tobii system, specifically
21 designed to engage young participants. This calibration was critical for accurate gaze position
22 tracking and was repeated as necessary, particularly in instances where the eye-tracking device
23 showed any discrepancies in detecting the participant's gaze. To maintain consistency and
24 reliability in data quality, we ensure constant lighting conditions in the testing room throughout all
25 sessions. Special consideration was given to the youngest participants, who were seated on their
26 parent's lap when they felt more comfortable in this setting, a strategy that effectively minimized
27 potential head and body movements that could interfere with accurate data collection. For data
28 analysis, we employed the Tobii IV-T Fixation filter (13,76). This tool is specifically designed to
29 extract fixation data, providing us with precise and reliable measures of visual attention and
30 engagement.

31

32 **2.2 Audio and visual stimuli**

33 We edited the original movie soundtrack using Audacity v.2.2.1 (Audacity Team, 2021) to isolate
34 speech excerpts while removing extraneous background noise such as birds' singing and musical
35 interludes.

36

37 Having extracted the video speech-track, we explored corresponding visual dynamics. Within the
38 video excerpts that contained speech, we distinguished two key components: speech envelope
39 (Figure 1B1) and visual motion (Figure 1B2). To extract the speech envelope, we used the absolute
40 value of the analytic signal (77). The obtained speech envelope was then down-sampled to 1000

1 Hz and filtered using a zero-phase, fourth-order Butterworth filter with a 40 Hz cutoff. The visual
2 motion was restricted to the participant's gaze-attended zone. We considered individual gaze-
3 fixation positions and the size of retinal stimuli around the gaze-fixation point. Our analysis took
4 into account the decline in visual acuity and the crowding effects of parafoveal vision (2–5 degrees
5 from the fixation point) (78). The region for capturing visual stimuli through eye gaze was defined
6 as a square with sides the length of an 8-degree diameter centered on the fixation point, aligning
7 with findings on the effective visual span guiding saccades (79). This corresponded to
8 approximately 318×318 pixels (Figure 1 A1 & A2). The gaze-capture zone was converted to
9 grayscale, and luminance differences between successive frames were computed following
10 established methodologies (80). The extraction of corresponding visual stimuli involved
11 converting the region of interest of each frame to a grayscale and computing the luminance
12 difference between successive frames. Pixels with a luminance change greater than 10 (a threshold
13 chosen to mitigate video recording noise) were selected, and the average luminance change
14 constructed the visual motion component.

15
16 Visual motion data was upsampled to match the 1000 Hz EEG sampling frequency. Speech
17 envelopes and visual motion corresponding to the same speech-track were aligned and prepared
18 for further analysis. With this method, visual motion is inherently individual-specific, as it relies
19 on the visual exploration of each participant.

20

21 **2.3 Stimulus features analysis**

22 In order to assess shared information between speech envelope and visual motion, in ASD and TD
23 groups, we calculated mutual information (MI) scores, a dynamic metric, expressed in bits, which
24 quantifies the reduction in uncertainty of one variable when another is observed (81,82). We
25 calculated MI using the *quickMI* function from the Neuroscience Information Theory Toolbox
26 (83). The parameters for this calculation were set to 4 bins, no delay, and a p-value threshold of
27 0.001 (83). For generating the MI scores, we concatenated all kept excerpts in the same sequence
28 across subjects, separately for each stimulus feature and each group (ASD and TD). This process
29 was followed by a comparative analysis of MI values between groups.

30

31 We only included stimuli corresponding to time periods with usable EEG signals. This resulted in
32 slight variations in the stimulus duration for the ASD and TD groups, for which we controlled. No
33 significant disparities in MI scores emerged between the two groups ($t(1, 61) = 1.250$, $p = 0.216$,
34 Cohen's $d = 0.315$, Figure 1C), indicating that these minor differences did not lead to notable
35 group differences in the shared information between the speech envelope and visual motion.

36

37 **3 EEG acquisition and pre-processing**

38 The EEG data were acquired through a 129-electrode (Hydrocel Geodesic Sensor Net (HCGSN)
39 system (Electrical Geodesics, USA) at a sampling rate of 1000 Hz. During recording, the signals
40 were subjected to real-time 0-100 Hz band-pass filtering. The reference electrode was positioned

1 at the vertex (Cz). Data pre-processing was conducted using the EEGLAB v2019 toolbox within
2 the MATLAB environment (84) and Cartool (<https://sites.google.com/site/cartoolcommunity/>).
3 One hundred and ten channels were kept excluding the cheek and neck to prevent contamination
4 by muscle artifacts. EEG signals were filtered using a zero-phase fourth-order Butterworth
5 bandpass (0.1-70 Hz) and a 50 Hz notch filter to eliminate power line noise. EEG data were
6 visually inspected to remove movement artifact-contaminated periods. Bad channels were first
7 identified and excluded for excessive signal amplitude. Eye blinks, saccades, electrical noise, and
8 heartbeat artifacts were excluded using independent component analysis (ICA). A spherical spline
9 interpolation was used to interpolate the channels contaminated by noise using the ICA-corrected
10 data. Finally, a common average reference was recalculated on the cleaned data, with an additional
11 step of applying a 30Hz low-pass filtering (80). To ensure that all the EEG signals and stimulus
12 features were on a similar scale and thus comparable, we normalized both the EEG signals and
13 stimulus features (i.e. speech envelope and visual motion) using the *nt_normcol* function
14 (Noisetools: <http://audition.ens.fr/adc/NoiseTools/>).

15

16 **4 Temporal Response Functions (TRF)**

17 To quantify how well EEG in ASD and TD children linearly varied with the stimulus features, we
18 performed regularized regression (with ridge parameter λ) as implemented in the mTRF toolbox
19 (85). The Temporal Response Function (TRF) captures how the brain's EEG activity correlates
20 with and responds to changes in the stimulus over time, providing a dynamic mapping of the neural
21 processing of the stimulus features. More precisely, the TRF accounts for the fact that the brain's
22 response to a stimulus occurs with a certain delay. Specifically, the TRF analysis models the
23 relationship between each stimulus feature (i.e. speech envelope or visual motion) and the brain's
24 response, particularly the time-lagged aspects of this relationship. The TRF includes a coefficient
25 for each time lag that quantifies the strength and direction (positive or negative) of the brain's
26 response to the stimulus at that specific time delay. The lags in a TRF are represented as a series
27 of time-points or intervals, typically in milliseconds. For the TRF estimation, we downsampled all
28 signals to a rate of 100 Hz to speed computation.

29

30 **4.1. Estimation of TRF using forward encoding models**

31 Specifically, we used a forward encoding model approach. Since changes in the EEG signal are
32 expected with an unknown time lag after the stimulus, predictions were computed over a range of
33 time lags between 300 ms earlier than the stimulus and 300 ms later than the stimulus. To single
34 out brain signals involved in auditory and visual dynamics, we constructed two distinct univariate
35 encoding models using the speech envelope and visual motion as independent regressors
36 respectively (labeled A-only and V-only). Further advancing our investigation, we developed a
37 multivariate encoding model, labeled 'AV-joint', which integrates both the speech envelope and
38 visual motion as regressors. The integration within this model is operationalized through the
39 assignment of trade-off weights to the AV regressor, which are calibrated to reflect the respective
40 contributions of auditory and visual stimuli. These trade-off weights ensure a balanced

1 representation within the model, enabling for an accurate representation of the brain's concurrent
2 processing of both modalities. This balanced approach aims to embody effective multisensory
3 integration, preventing the overshadowing of one sensory modality by another.

4
5 The comparative analysis of the AV-joint model against the A-only and V-only models is essential
6 to elucidate the incremental benefits of simultaneous audio-visual integration over unimodal
7 processing. By evaluating the performance and predictive accuracy of the multivariate model
8 relative to its univariate counterparts (one focused only on auditory processing and the other on
9 visual processing), we aim to substantiate the hypothesis that the synergistic consideration of
10 auditory and visual stimuli offers a more comprehensive understanding of sensory processing in
11 naturalistic listening environments.

12
13 For each group, we created "generic" models that predict the EEG data of an individual participant
14 (n th participant) using a TRF derived from the EEG data of the other participants (the remaining
15 $n-1$ participants). We thus implemented an n -fold leave-one-out cross-validation strategy and
16 optimized the model through a parameter search for the regularization parameter λ .

17

18 **4.2 Selection of optimal regularization parameter**

19 To mitigate the risk of data overfitting in forward encoding models, we integrated an optimized
20 regularization parameter λ , determined for each stimulus feature. To do so, we trained multiple
21 model iterations on subsets of data comprising $n-1$ participants. During these iterations, λ was
22 systematically adjusted within a predefined range from 10^0 to 10^5 , with increments in the exponent
23 of 0.5. The criterion for selecting the optimal λ was maximal predictive accuracy. This was
24 quantitatively assessed by calculating Pearson's correlation coefficient (r) between the predicted
25 EEG signals and the observed signals, for each electrode and each participant. Once optimal λ was
26 identified, the refined model was used to predict the EEG responses of the n th participant, using
27 an n -fold leave-one-out cross-validation paradigm. This methodology was applied consistently
28 across both ASD and TD groups, ensuring the robustness and reliability of the predictive models
29 within each group.

30

31 **4.3. Pearson's correlation coefficient (r)**

32 Pearson's correlation coefficient (r) was computed to quantify the model's prediction accuracy. For
33 each electrode, the correlation between the EEG signals and the TRF-predicted signals was
34 calculated across all time points. This process was repeated for each participant for the validation
35 set. The correlation coefficients were then averaged across participants to obtain a mean Pearson's
36 r -value per electrode. This allowed us to create a scalp-wide map of the model's prediction
37 accuracy, highlighting areas with the strongest correlation between the predicted and observed
38 EEG responses. The accuracy of EEG predictions derived from the optimized model was then
39 converted into Z -scores for further statistical group comparisons. The Z -score computation
40 involved subtracting the mean of surrogate data and dividing it by the standard deviation of

1 surrogate data. Surrogate distributions were generated by randomly shifting (50 times) the orders
2 of the testing EEG segments, preserving their original temporal structure. In summary, this
3 analysis indicated the accuracy with which stimulus features are predicted from the EEG data for
4 each participant.

5
6 To quantitatively compare the accuracy between the ASD and TD groups, we used a cluster-based
7 permutation test with 1000 randomization iterations, following the approach of Maris and
8 Oostenveld (86). Clusters were defined by considering both time and spatial electrode
9 configurations, requiring each cluster to include at least two adjacent electrodes. A pivotal aspect
10 of this approach was ensuring that the cluster-level type-I-error probability remained below the
11 0.05 threshold. This strategy was effective in controlling the family-wise error rate, maintaining it
12 within the 5% type-I-error rate boundary.

13

14 **4.4. Stimulus reconstruction using decoding models**

15 We trained decoders using EEG data across a wide temporal range, from -300 to 300 ms relative
16 to the stimulus, aiming for optimal stimulus reconstruction. This analysis involved an iterative
17 process of leave-one-out cross-validation and optimization techniques to refine our decoding
18 accuracy, assessing how effectively the EEG signals could predict the stimulus features (i.e., visual
19 motion and speech envelope) for each group.

20

21 To enhance our understanding of the temporal alignment between EEG signals and the stimuli,
22 our methodology evolved to focus on discrete, predefined time-lag intervals for decoder training,
23 rather than a continuous range. By pinpointing discrete intervals and evaluating their decoding
24 success, we identified the most informative time-lag segment that yielded the highest fidelity in
25 stimulus reconstruction. Determining this optimal time-lag interval informs us about the specific
26 moments when the EEG data are most synchronized with the stimulus features, thus shedding light
27 on the precise neural timing critical for effective sensory processing and integration.

28 For the decoding models, we determined the accuracy of stimulus reconstruction accuracy and the
29 identity of the best time lag using the Kruskal-Wallis test with Dunn's multiple comparisons test
30 (87). This non-parametric statistical method was chosen for its capability to handle variations in
31 group means and variances across different conditions.

32

33 **5. Low-frequency tracking of audio-visual signal**

34 To explore whether the combined processing of auditory and visual stimuli can be explained by
35 the tracking of audio-visual signals by low-frequency brain activity, we used a coherence-based
36 and phase-based analytical framework (88). This approach probed the interplay between neural
37 responses and stimulus features by comparing their magnitude spectra and the phase relationship.
38 Our focus centered on the delta and theta frequency bands, which are critical for the temporal
39 coordination necessary for effective integration of multimodal information (89).

40

1 **5.1. Coherence Analysis**

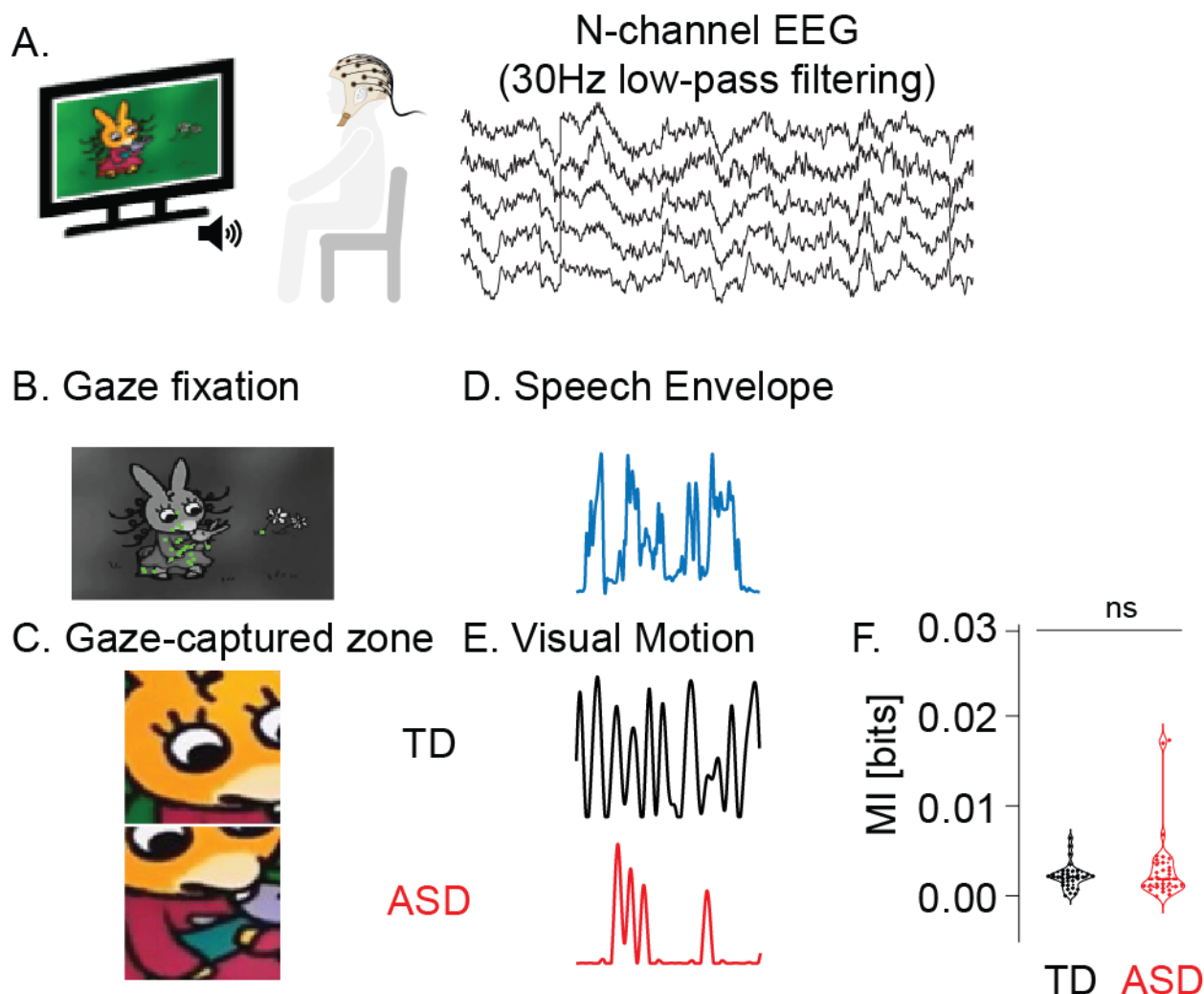
2 We assessed individual responses to speech envelope and visual motion by computing magnitude-
3 squared coherence for each trial and electrode using the *mcoh* function in Matlab, applying
4 Welch's averaged modified periodogram method. The analysis spanned a frequency range from
5 0.1 to 30 Hz, in increments of 0.33 Hz (90).

6
7 The analysis targeted delta (δ , ~4 Hz) and theta (θ , 4-8 Hz) frequency bands, identifying
8 frequencies where coherence peaked most prominently for each stimulus condition. Statistical
9 comparisons across groups and stimuli were conducted using the clusters identified through the
10 method outlined in Section 4, combined with a nonparametric test and Dunn's multiple
11 comparisons test (87). Significance of these findings was established using a surrogate-corrected
12 coherence approach. Surrogate distributions were generated by randomly shifting the neural time
13 course relative to the stimulus feature time courses, preserving their original temporal structure.
14 This process was repeated 50 times for each stimulus condition to generate a robust surrogate
15 distribution. The resulting coherence values were then standardized (Z-scored) against this
16 distribution.

17 **5.2. Phase Analysis**

18 We also performed a phase analysis by calculating the cross-power spectral density (CPSD) phase
19 for each stimulus, electrode, and trial. This was done using the *cpsd* function in Matlab, employing
20 parameters consistent with the coherence analyses. Phase values were determined based on the
21 peak frequency identified in the coherence analysis. Group comparisons were conducted using
22 Matlab's Circular Statistics Toolbox (91). A two-way parametric ANOVA for circular data was
23 performed to facilitate a nuanced comparison between pairs of conditions and groups, followed by
24 post-hoc comparisons using the Watson-Williams multi-sample test (91). Meanwhile, the Rayleigh
25 test was performed to investigate whether phase distribution is under unimodal distribution. Our
26 focus was primarily on electrodes identified through TRF estimation outcomes.

27
28



1
2 **Figure 1. Overview of Experimental Procedures and Features of Interest.** (A) Experimental procedures (B) Gaze
3 fixation. Example of individual gaze fixation points (green dots) on a black and white image; (C) Example of gaze-
4 captured screen areas. Depiction of screen areas captured by the gaze of participants in ASD (Autism Spectrum
5 Disorder) and TD (Typically Developing) groups. (D) Example of a stimulus speech envelope from the cartoon
6 soundtrack. (E) Visual motion corresponds to the same stimulus in each group. (F) Comparison of speech envelope
7 and visual motion. Mutual Information (MI) between ASD and TD groups (ns. $p>0.05$).
8

9 Results

10 1.1 Atypical neural tracking of speech envelope in ASD children

11 We found distinct neural tracking for the auditory and visual parts of the stimuli. Different scalp
12 distribution patterns were observed between the two groups (Figure 2). The ASD group had
13 reduced neural response to the speech envelope relative to the TD group (Figure 2A, top row). In
14 contrast, there was no difference in visual motion processing between the two groups (Figure 2A
15 middle row). These results suggest atypical auditory but not visual processing of dynamic
16 communicative stimuli in young children diagnosed with ASD.
17

1 1.2 Univariate decoding models confirm atypical speech envelope processing in children with 2 ASD

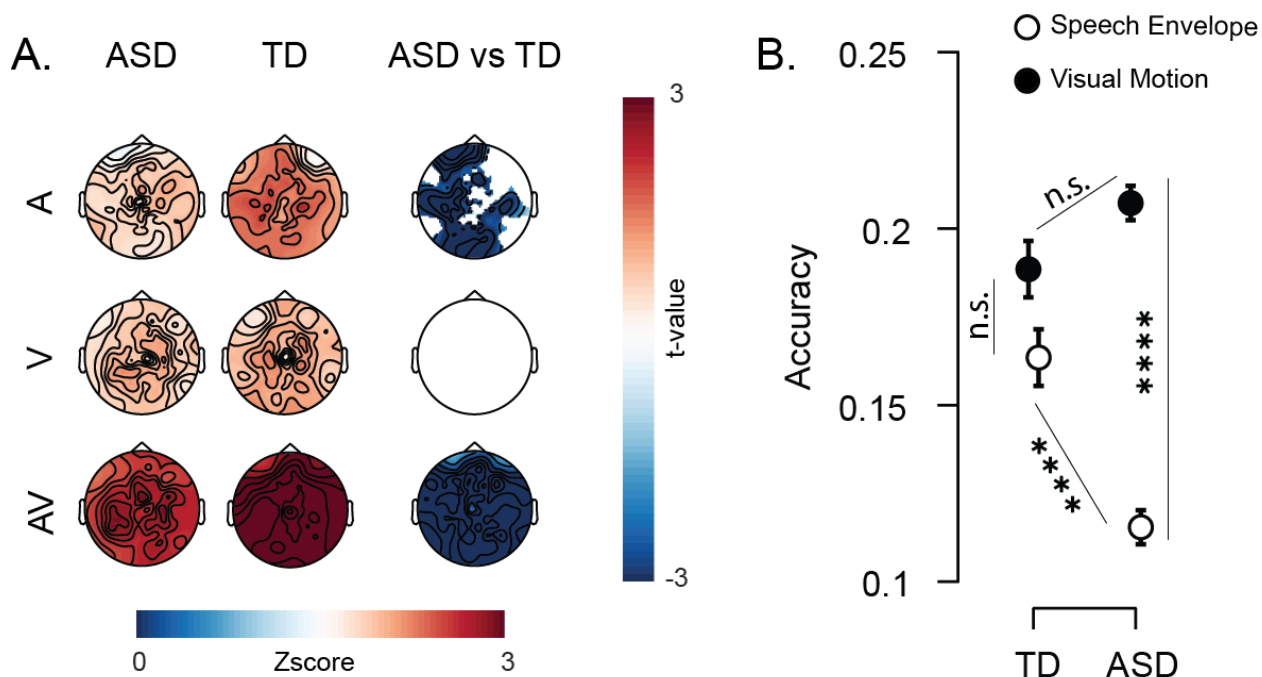
3 Stimulus reconstruction accuracy was different in the ASD and TD groups. While reconstruction
4 accuracy was equivalent for speech envelope and visual motion ($p > 0.9999$, Figure 2B) in the TD
5 group, it was lower for speech than for visual motion ($p < 0.0001$, Figure 2B) in the ASD group.
6 Group comparison shows that speech envelope reconstruction was weaker in the ASD than in the
7 TD group ($p < 0.0001$, Figure 2B), without significant group difference for visual motion
8 reconstruction (Figure 2B). These results align with neural tracking results to suggest that it is
9 mostly speech processing that is impaired in autistic children.

10

11 2. Multivariate encoding models indicate atypical audiovisual processing in ASD children

12 Further, we explored whether speech anomalies in ASD are limited to auditory processing
13 difficulties or associated with audiovisual (AV) processing anomalies. We found stronger neural
14 representation of the combined AV stimulus compared to individual single stimuli in both groups,
15 an expected result as multivariate models provide in general more accurate prediction of neural
16 responses by combining information from multiple sources. More interestingly, the joint model
17 suggested weaker AV representation in the ASD than in the TD group (Figure 2A bottom row).
18

18



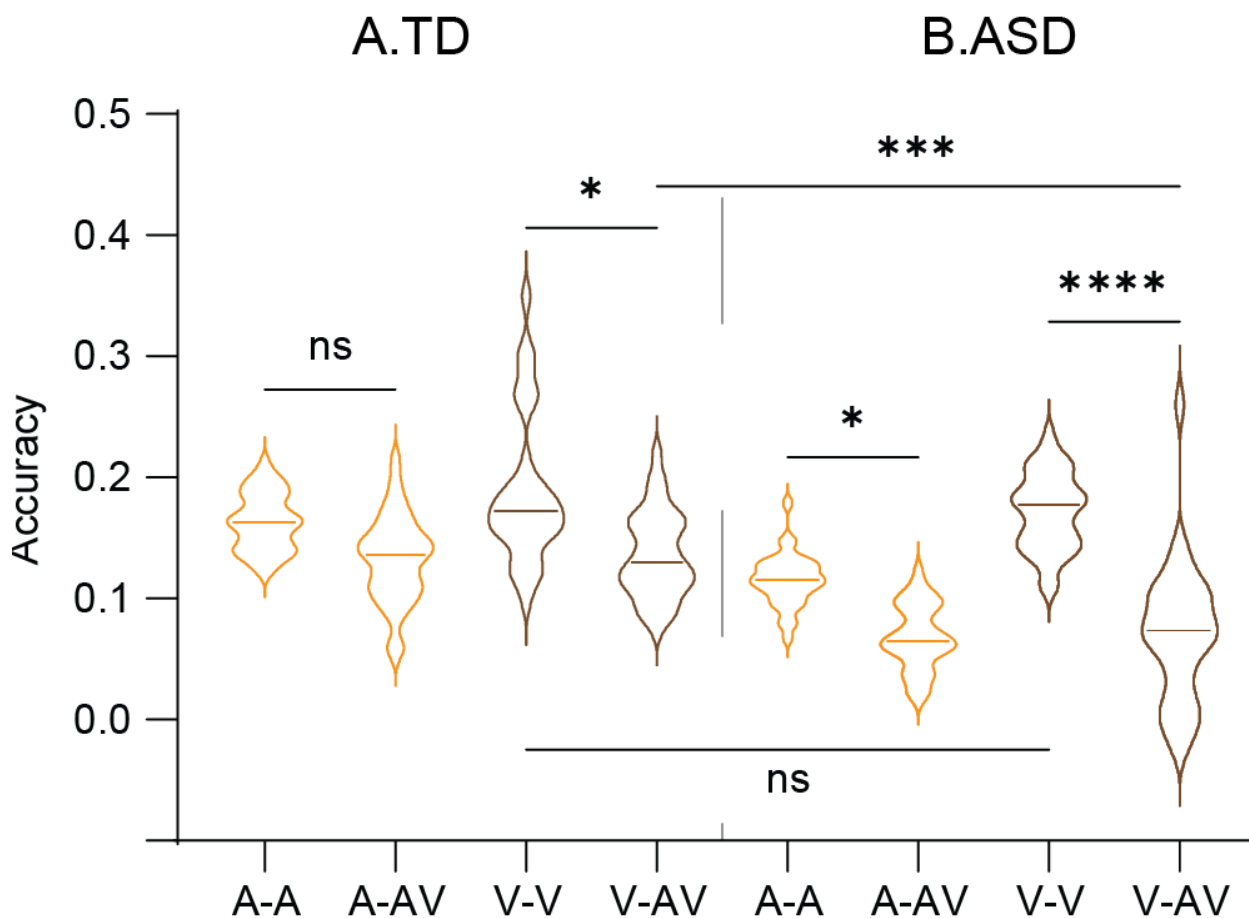
19

20 **Figure 2. Comparison of audio, visual, and AV models.** (A.) *Neural representations* in ASD (left) and TD (middle)
21 groups, for each model (A-speech envelope, V-visual motion, and AV joint) across all scalp electrodes. The right
22 column shows EEG channels where significant group differences are observed using cluster-based nonparametric
23 statistics ($p < 0.05$; with a positive t -value indicating greater predictability in the ASD group compared to the TD
24 group). (B.) **Stimulus reconstruction accuracy** for speech envelope and visual motion in both groups. Error bars
25 indicate the standard error of the mean. Significance levels are indicated as follows: 'ns' for $p > 0.05$ (not significant),
26 * for $p < 0.05$, ** for $p < 0.01$, *** for $p < 0.001$, and **** for $p < 0.0001$.

1
2 **3. Compromised audiovisual integration in ASD disrupts visual enhancement of auditory**
3 **processing**

4 By comparing the decoding accuracies for speech envelope and visual motion across univariate
5 (A-only and V-only) and multivariate (AV-joint) models, we found distinct patterns of audiovisual
6 integration in ASD and TD children (Table 1, Figure 3). Notably, in the TD group, the accuracy
7 of speech envelope reconstruction in the AV-joint model (concurrent speech envelope and visual
8 motion) did not significantly differ from the A-only model ($p = 0.1414$, Figure 3A). Conversely,
9 in the ASD group, the speech envelope was decoded with significantly less accuracy in the AV-
10 joint model than in the A-only model ($p=0.0304$, Figure 3B), indicating a disruptive effect of AV
11 integration on auditory processing specific to this group. Despite these differences, both groups
12 exhibited decreased visual motion reconstruction accuracy in the AV-joint model compared to the
13 V-only model (ASD group: $p < 0.0001$, TD group: $p = 0.0373$, Figure 3), with a more pronounced
14 decrement observed in the ASD group ($p = 0.0003$, Table 2). These findings suggest that while
15 AV integration generally impacts visual processing across both groups, it adversely affects
16 auditory processing primarily in the ASD group.

17
18



19
20 *Figure 3. Evaluation of decoding accuracy in ASD (A) and TD (B). Stimulus reconstruction accuracy: speech*

1 *envelope(A-) and visual motion(V-) in both the single-stimulus model (A-only = A-A and V-only = V-V) and the AV-*
2 *joint model (A-AV, and V-AV)). Significance levels are indicated as follows: 'ns' for $p>0.05$ (not significant), * for p*
3 *<0.05 , ** for $p<0.01$, *** for $p<0.001$, **** for $p<0.0001$. For additional details, see Supplemental Figure 3-1.*

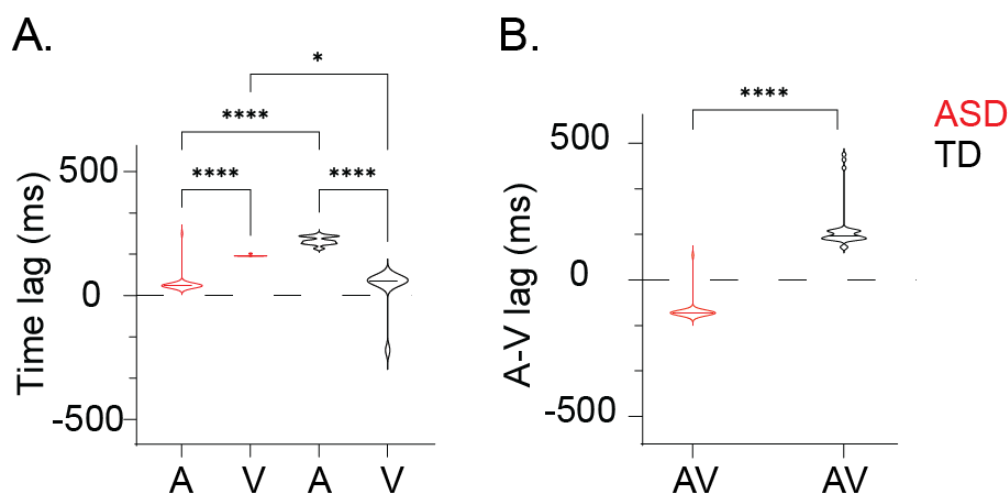
4

5 **4. No visual precedence in audiovisual processing in ASD**

6 When studying the temporal dynamics of auditory and visual processing to assess audiovisual
7 integration, we found distinct time-lags in auditory and visual decoding accuracies in ASD and
8 TD children. In the TD group, it took 200~ms to reach significant decoding accuracy for auditory
9 responses but only ~50ms for visual ones. The exact opposite pattern was found in the ASD group
10 with a ~200 ms time-lag got visual decoding versus ~50 ms for auditory decoding (Figure 4A).
11 This analysis revealed a visual lead in the TD consistent with the fact that speech sources are
12 usually ahead of sounds, but an auditory lead in the ASD group (Figure 4B). This temporal
13 inversion indicates a fundamental alteration in the sensory processing sequence for children with
14 ASD.

15

16



17

18

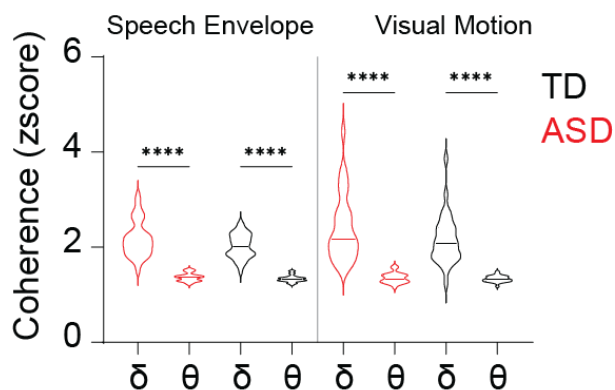
19 **Figure 4. Optimal EEG-stimuli time lag for ASD (red) and TD (black) groups.** (A) depicts the optimal time-lag
20 observed in the reconstruction of stimulus features in AV-joint model, specifically speech envelope (A-) and visual
21 motion(V-); Positive values represent stimulus lead EEG signal. (B) illustrates the A-V time lag AV-joint model.
22 Positive values represent V leads A. Significance levels are indicated as follows: ns>0.05, * $p <0.05$, ** $p<0.01$,
23 *** $p<0.001$, **** $p<0.0001$.

24

25 **5. Intact capacity to synchronize neural activity to the stimulus input in ASD**

26 To ascertain whether speech tracking anomalies in ASD are primarily attributable to a general
27 defect in stimulus/brain synchronization or rather results from audiovisual integration deficits, we
28 investigated the intricate temporal dynamics underlying these processes. Specifically, we explored
29 the coherence of stimulus-response relationships within delta (1–4 Hz) and theta (4–8 Hz)
30 frequency bands, typically associated with syllable-and phrase-level speech processing. We
31 observed a higher stimulus/brain coherence in the delta band than the theta band, yet when

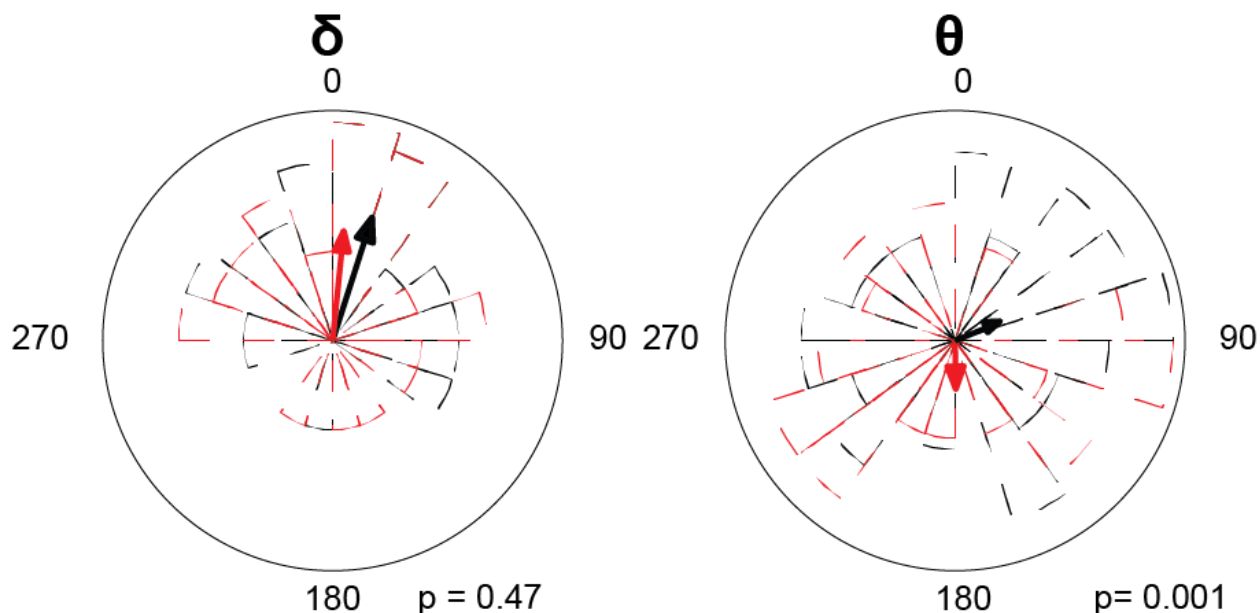
1 comparing across groups, both the average coherence values and their distribution patterns showed
2 remarkable similarity (Figure 5, Table 3). This uniformity shows that the intrinsic capacity of
3 neural activity to synchronize with auditory and visual stimuli is consistent between groups (ASD
4 v.s. TD).
5



6
7 **Figure 5. Stimulus-response coherence** in Theta and Delta Bands for ASD (red) and TD (black) groups. The plot
8 displays the coherence between stimulus and response for Speech Envelope and Visual Motion. Error bars represent
9 the standard error of the mean. The coherence levels are compared within the specific frequency bands of interest,
10 highlighting potential group differences in sensory processing. Significance levels are indicated as follows: ns > 0.05,
11 *p < 0.05, **p < 0.01, ***p < 0.001, ****p < 0.0001.

12 **6. Theta-range desynchronization of audio-visual responses in ASD**

13 Given the preserved synchronization capacity between brain activity and external stimuli in each
14 modality, we then sought whether audio-visual integration anomalies, notably the inverted AV
15 temporal patterns, are associated with a phase desynchronization of auditory and visual processing.
16 In the delta band, we observed similar phase angles for both groups ($F(1,62) = 0.494$, $p < 0.470$),
17 indicating comparable phase locking at this slower frequency, suggesting that temporal alignment
18 in this low-frequency range does not differentiate between groups. Moreover, the small angle
19 indicates the absence of delta band phase-shift between modalities. Yet, a significant group
20 difference was observed in the theta band. The TD group showed a phase-shift of approximately
21 90 degrees, signaling effective sequential integration with one modality leading the other by a
22 consistent temporal offset, that optimizes audio-visual integration at the syllable level. In ASD
23 children the phase-shift amounted to 180 degrees and the group difference was significant ($F(1,62)$
24 $= 12.05$, $p < 0.001$) (Figure 6). The observed 180-degree phase shift in ASD could suggest that
25 auditory and visual information is out of sync: when one sensory modality is at its peak processing
26 efficiency, the other is at its lowest, potentially leading to disjointed even conflicting sensory
27 processes.



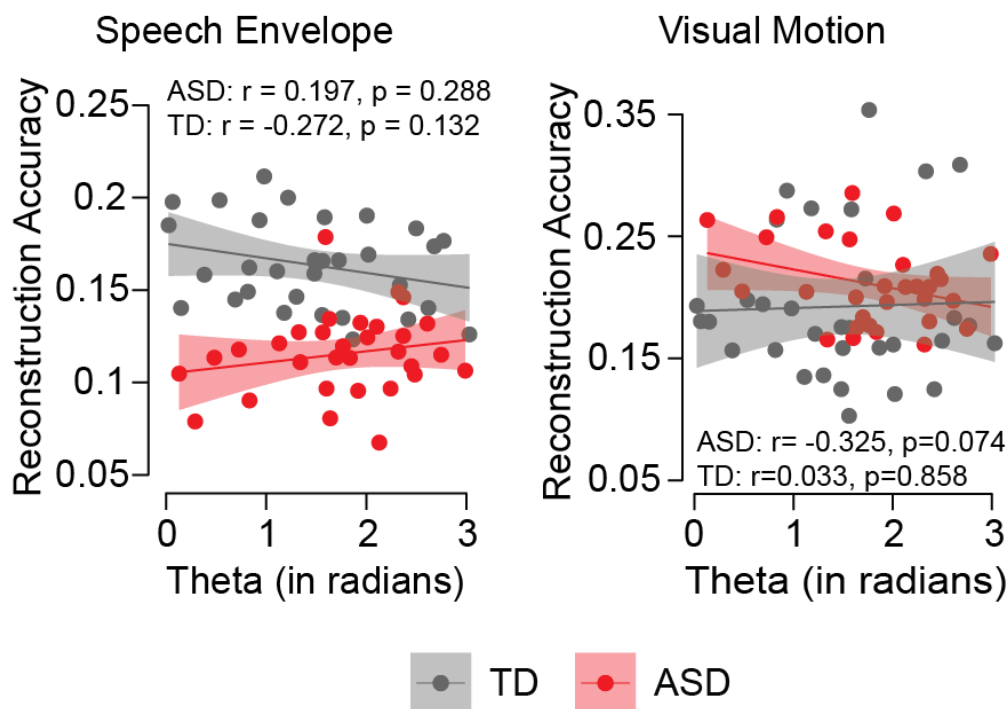
1
2 **Figure 6. Phase-shift distribution between speech envelope and visual motion.** This figure shows the phase-shift
3 distribution between the brain processes of speech envelope and visual motion stimuli for each group. The circular
4 mean of the phase-shift across all subjects is indicated by colored lines: red for the ASD group and black for the TD
5 group. Corresponding polar histograms in red (ASD) and black (TD) visually represent the distribution of phase-
6 shifts for each group. Both groups were tested against the hypothetical uniform distribution of delta (rayleigh test,
7 ASD: $p < 0.001$, rayleigh $r = 0.98$, TD: $p < 0.001$, rayleigh $r = 0.98$) and theta phase (rayleigh test, ASD: $p < 0.001$,
8 rayleigh $r = 0.95$, TD: $p < 0.001$, rayleigh $r = 0.96$).

9
10 **7. AV phase-shift is related to auditory encoding accuracy in TD and visual encoding**
11 **accuracy in ASD**

12 Finally, we sought to understand the relationship between the theta band AV phase-shift and the
13 accuracy of auditory and visual information reconstruction within an unimodal framework (Figure
14 7). As could be expected, in TD children the AV phase-shift did not influence visual reconstruction
15 accuracy ($r = 0.033$, $p = 0.858$), but there was a weak negative correlation between the phase-shift
16 extent and speech reconstruction accuracy ($r = -0.272$, $p = 0.132$): when the AV phase-shift
17 increased speech reconstruction accuracy decreased, which given the visual lead previously
18 observed could suggest a causal effect. A different pattern was seen in ASD children, with no
19 relation between the phase shift extent and speech reconstruction accuracy ($r = 0.197$, $p = 0.288$),
20 but a weak negative correlation between the phase-shift extent and the accuracy of visual
21 information reconstruction ($r = -0.325$, $p = 0.074$), with larger AV phase-shifts linked to poorer
22 visual reconstruction accuracy. Likewise, given the auditory lead observed in children with ASD,
23 this could suggest a causal effect.

24 Contrasting patterns in TD and ASD children underscore distinct audiovisual integration
25 mechanisms. In logic, reconstruction accuracy is a proxy of sensory encoding accuracy. Thus, in
26 TD children, the data confirm the known reliance on visual cues to enhance auditory processing,
27 with any misalignments adversely affecting speech information integration. Conversely, in ASD

1 children, while speech encoding is weaker and overall less dependent on visual-auditory phase
2 congruency, visual processing is vulnerable to strong AV desynchronization.
3



4
5 **Figure 7. The relationship between theta phase-shift and reconstruction accuracy of speech envelope (left) and**
6 **visual motion (right) in ASD and TD. ASD group suggests a greater phase shift between speech envelope and**
7 **visual motion positively correlates with speech reconstruction accuracy but negatively correlates with visual reconstruction,**
8 **while** reversely **in** TD **group.**
9

10 Discussion

11 Using several analyses of the EEG recorded in very young children with and without ASD while
12 they were watching a short animated movie, we confirmed previous results showing profound
13 anomalies of the capacity to follow speech rhythms (24,92), an essential prerequisite to speech
14 comprehension. The present study goes beyond this observation by showing that children with
15 ASD did not exhibit the natural dominance of auditory processing when exposed to natural audio-
16 visual speech conditions. Instead, their processing of audio-visual stimuli was impacted by a
17 temporal misalignment of these sensory inputs, which disturbs the predictive processing typically
18 at play when perceiving speech.
19

20 Audio-visual integration anomalies interfere with sensory encoding in ASD

21 Audio-visual processing, especially the synchronization of the two sensory modalities, plays a
22 pivotal role in understanding the communicative challenges observed in ASD. Previous studies
23 have established that basic auditory dysfunctions and atypical speech processing are characteristic
24 of ASD from an early developmental stage (19–24). Our study reveals that such anomalies are not
25 merely isolated auditory deficits but are deeply connected to the integration of auditory and visual

1 information, a process critical for effective communication, particularly in dynamic or complex
2 listening environments (93,94).

3
4 Our findings reveal a specific disruption in audio-visual integration among children with ASD,
5 manifesting in visual dominance and temporal disorganization in auditory and visual processing.
6 This disruption sharply contrasts with the expected auditory processing dominance (95) and might
7 significantly contribute to the language development anomalies encountered by these children. In
8 typical development, the precedence of orofacial visual cues during speech facilitates auditory
9 comprehension through predictive processing, optimizing the brain's synchronization to incoming
10 speech signals (96). In ASD, extended integration time windows and the lack of effective
11 synchronization of auditory responses by visual signals (as evidenced by the atypical theta band
12 phase-shifts) suggest they cannot use visual cues to facilitate auditory speech processing. On the
13 contrary, our findings show that auditory cues perturb visual processing of communicative
14 situations.

15

16 **Specific repercussions of disrupted audio-visual processing on speech tracking**

17 Children with ASD demonstrate effective visual motion tracking and processing capabilities,
18 comparable to their TD peers, despite distinctive scene analysis patterns as previously observed
19 (13,97). Within their preferred exploration zones, children with ASD process visual motion
20 similarly to their TD peers (97), exhibiting a level of bottom-up excitability to visual stimuli akin
21 to TD children (98,99). Our univariate encoding results suggest that the neural activity responsible
22 for visual motion tracking operates similarly in both ASD and TD groups.

23
24 However, when visual processing co occurs with speech processing, some difficulties appear. Our
25 multivariate modeling indicates that the neural encoding of audio-visual percepts in ASD children
26 is less efficient, confirming that audiovisual contexts can disrupt brain responses to speech in this
27 population (100). In the same vein, Shic et al. (2020) underscore such AV integration difficulties,
28 noting that children with ASD tend to attend less to faces and mouths in general and more
29 specifically when they produce speech (7). Our study reinforces this crucial finding by showing
30 that while children with ASD encode single visual streams relatively well (visual motion tracking
31 in univariate modeling), they struggle with the concurrent encoding of both auditory and visual
32 streams (multivariate modeling).

33
34 Building upon this framework, the research conducted by Chawarska et al. (2022) reveals that 12-
35 month-old infants at risk for ASD, even though they explore faces and mouths similarly to infants
36 with no family history of autism (101), cannot leverage audiovisual cues for language acquisition
37 as do typical children. Our study uncovers the potential underpinnings of the audiovisual
38 integration difficulties observed in ASD. The decoding results indicate that while audio-visual
39 integration interferes with visual processing in both ASD and TD groups, its impact on speech
40 processing is particularly detrimental in the ASD group. Thus, the impairments in AV integration

1 we observe are not merely additive but they interactively exacerbate sensory processing challenges
2 much more adversely in children with ASD.

3

4 **Audio-visual temporal integration underlies speech impairment in ASD**

5 Audiovisual integration capitalizes on the temporal alignment of sensory events, with visual
6 information often enhancing the auditory signal's clarity and precision, especially when auditory
7 cues are poor, noisy, or ambiguous (102–104). Visual cues related to speech are typically
8 processed faster than auditory cues, allowing visual information to facilitate synchronizing
9 subsequent auditory processing (105). Our findings confirm in TD children a visual lead (~50 ms)
10 within a temporal window is conducive to effective interaction and coordination between auditory
11 and visual cues. This window reasonably aligns with established models, positing a 200 ms
12 integration period (39,61–63), ranging from a 30 ms visual lag to a 170 ms of visual lead (61).

13

14 The precise timing of audio-visual sequences is fundamental to audio-visual integration via
15 predictive processing, whereby the brain leverages visual cues to anticipate and decode
16 forthcoming auditory information. Here, phase-locking analyses in TD children show that the
17 neural responses associated with auditory and visual processing exhibit a 90-degree phase shift.
18 This observation indicates that the brain orchestrates visual and auditory information in a
19 synergistic but temporally distinct manner. Such a phase relationship is instrumental in achieving
20 a dynamic balance between the sensory streams, facilitating an integration that enhances
21 perception and communication (106,107). The pivotal role of the theta frequency band in
22 orchestrating audio-visual speech processing is robustly supported in the literature (56,58,59). A
23 $\pi/2$ phase shift during an audio-visual speech event might fine-tune the phase alignment to a
24 timing that is congruent with the auditory inputs. This meticulous phase alignment contrasts
25 sharply with the broad phase distribution observed in the ASD group, hinting at a pronounced
26 disparity in how auditory and visual cues are integrated.

27

28 Crucially, in ASD children, the integration of auditory and visual streams is jeopardized, as
29 evidenced by our observation of an atypical auditory lead (~50 ms), which disrupts the
30 conventional sequence where visual information typically precedes auditory. This inversion
31 undermines the usual enhancement provided by visual signals to auditory processing, highlighting
32 a marked alteration or impairment in multisensory integration. Furthermore, we noted a 180-
33 degree phase-shift in the neural activities associated with processing these 2 streams. This
34 substantial phase misalignment reflects a profound disruption in temporal coordination, potentially
35 leading to confusion or interpretation errors. Such a discrepancy underscores a critical deficiency
36 in predictive processing in ASD, where, rather than synergistically enhancing each other, auditory
37 and visual cues conflict, undermining the synthesis of coherent audio-visual perception. This
38 misalignment is also reflected in the broader phase distribution seen in children with ASD,
39 suggesting that they might require an extended temporal window to reach effective processing
40 (31,32).

1
2 Our results thus confirm that the phase of low-frequency neural oscillations is crucial for the
3 encoding of order - for instance with the implication of the theta band in working memory (108)
4 or for temporal parsing in speech (109). The anomaly in temporal encoding mechanisms described
5 in our experiment is constrained by the temporal features provided by external stimulation to build
6 a temporal reference frame. While delta oscillations have previously been linked to temporal
7 predictability (110,111), we observed here that sensory integration is affected by AV misalignment
8 in the theta range, which is associated with atypical speech perception in ASD. The AV integration
9 primarily occurs at the syllable level with a typical tolerance of AV asynchrony at 250ms, which
10 corresponds to the theta range (39,61–64).

11

12 **Conclusion**

13 We show remarkable anomalies in audio-visual integration in children with ASD. We confirm
14 previous findings of disrupted speech rhythm tracking and further reveal a specific disruption in
15 audio-visual integration, manifesting as temporal desynchronization. This disruption significantly
16 impacts speech processing, contributing to the communicative challenges faced by children with
17 ASD. Our results highlight the critical role of temporal processing in audio-visual integration and
18 underscore the importance of characterizing these mechanisms in ASD. Moving forward, these
19 insights could inform the development of targeted interventions aiming at regulating temporal
20 speech processing and AV synchronization to improve communication in children with ASD.

21

22 **References**

- 23 1. Alegria J, Noirot E (1978): Neonate orientation behaviour towards human voice.
24 *Int J Behav Dev* 1: 291–312.
- 25 2. Jusczyk PW, Bertoncini J (1988): Viewing the development of speech
26 perception as an innately guided learning process. *Lang Speech* 31: 217–
27 238.
- 28 3. Lecanuet J, Granier-Deferre C, Decasper A, Maugeais R, Andrieu A, Busnel M
29 (1987): Fetal perception and discrimination of speech stimuli; demonstration
30 by cardiac reactivity; preliminary results. *Comptes Rendus Académie Sci Sér*
31 *III Sci Vie* 305: 161–4.
- 32 4. Williams JHG, Massaro DW, Peel NJ, Bosseler A, Suddendorf T (2004):
33 Visual–auditory integration during speech imitation in autism. *Res Dev*
34 *Disabil* 25: 559–575.
- 35 5. Iarocci G, Rombough A, Yager J, Weeks DJ, Chua R (2010): Visual influences
36 on speech perception in children with autism. *Autism* 14: 305–20.
- 37 6. Feng S, Wang Q, Hu Y, Lu H, Li T, Song C, *et al.* (2023): Increasing

- 1 audiovisual speech integration in autism through enhanced attention to
2 mouth. *Dev Sci* 26: e13348.
- 3 7. Shic F, Wang Q, Macari SL, Chawarska K (2020): The role of limited salience
4 of speech in selective attention to faces in toddlers with autism spectrum
5 disorders. *J Child Psychol Psychiatry* 61: 459–469.
- 6 8. Klin A (1991): Young autistic children’s listening preferences in regard to
7 speech: a possible characterization of the symptom of social withdrawal. *J*
8 *Autism Dev Disord* 21: 29–42.
- 9 9. Klin A, Lin DJ, Gorrindo P, Ramsay G, Jones W (2009): Two-year-olds with
10 autism orient to non-social contingencies rather than biological motion [no.
11 7244]. *Nature* 459: 257–261.
- 12 10. Klin A (1992): Listening preferences in regard to speech in four children with
13 developmental disabilities. *J Child Psychol Psychiatry* 33: 763–9.
- 14 11. Dawson G, Meltzoff AN, Osterling J, Rinaldi J, Brown E (1998): Children
15 with autism fail to orient to naturally occurring social stimuli. *J Autism Dev*
16 *Disord* 28: 479–85.
- 17 12. Kuhl PK, Coffey-Corina S, Padden D, Dawson G (2005): Links between social
18 and linguistic processing of speech in preschool children with autism:
19 behavioral and electrophysiological measures. *Dev Sci* 8: F1–F12.
- 20 13. Kojovic N, Cekic S, Castañón SH, Franchini M, Sperdin HF, Sandini C, *et al.*
21 (2024): Unraveling the developmental dynamic of visual exploration of
22 social interactions in autism ((C. Büchel, editor)). *eLife* 13: e85623.
- 23 14. Federici A, Parma V, Vicovaro M, Radassao L, Casartelli L, Ronconi L (2020):
24 Anomalous Perception of Biological Motion in Autism: A Conceptual
25 Review and Meta-Analysis [no. 1]. *Sci Rep* 10: 4576.
- 26 15. Knight EJ, Krakowski AI, Freedman EG, Butler JS, Molholm S, Foxe JJ
27 (2022): Attentional influences on neural processing of biological motion in
28 typically developing children and those on the autism spectrum. *Mol Autism*
29 13: 33.
- 30 16. Todorova GK, Hatton REM, Pollick FE (2019): Biological motion perception
31 in autism spectrum disorder: a meta-analysis. *Mol Autism* 10: 49.
- 32 17. Marco EJ, Hinkley LB, Hill SS, Nagarajan SS (2011): Sensory processing in
33 autism: a review of neurophysiologic findings. *Pediatr Res* 69: 48R-54R.
- 34 18. Todorova GK, Pollick FE, Muckli L (2021): Special treatment of prediction
35 errors in autism spectrum disorder. *Neuropsychologia* 163: 108070.

- 1 19. Collet L, Roge B, Descouens D, Moron P, Duverdy F, Urgell H (1993):
2 Objective auditory dysfunction in infantile autism. *The Lancet* 342: 923–
3 924.
- 4 20. Haesen B, Boets B, Wagemans J (2011): A review of behavioural and
5 electrophysiological studies on auditory processing and speech perception in
6 autism spectrum disorders. *Res Autism Spectr Disord* 5: 701–714.
- 7 21. Edgar JC, Fisk Iv CL, Berman JI, Chudnovskaya D, Liu S, Pandey J, *et al.*
8 (2015): Auditory encoding abnormalities in children with autism spectrum
9 disorder suggest delayed development of auditory cortex. *Mol Autism* 6: 69.
- 10 22. Foss-Feig JH, Schauder KB, Key AP, Wallace MT, Stone WL (2017):
11 Audition-specific temporal processing deficits associated with language
12 function in children with autism spectrum disorder. *Autism Res Off J Int Soc*
13 *Autism Res* 10: 1845–1856.
- 14 23. Wang X, Wang S, Fan Y, Huang D, Zhang Y (2017): Speech-specific
15 categorical perception deficit in autism: An Event-Related Potential study of
16 lexical tone processing in Mandarin-speaking children. *Sci Rep.*
17 <https://doi.org/10.1038/srep43254>
- 18 24. Wang X, Delgado J, Marchesotti S, Kojovic N, Sperdin HF, Rihs TA, *et al.*
19 (2023): Speech Reception in Young Children with Autism Is Selectively
20 Indexed by a Neural Oscillation Coupling Anomaly. *J Neurosci* 43: 6779–
21 6795.
- 22 25. Benasich AA, Gou Z, Choudhury N, Harris KD (2008): Early cognitive and
23 language skills are linked to resting frontal gamma power across the first 3
24 years. *Behav Brain Res* 195: 215–22.
- 25 26. Benitez-Burraco A, Murphy E (2016): The Oscillopathic Nature of Language
26 Deficits in Autism: From Genes to Language Evolution. *Front Hum*
27 *Neurosci* 10: 120.
- 28 27. Cermak CA, Arshinoff S, Ribeiro de Oliveira L, Tendra A, Beal DS, Brian J,
29 *et al.* (2022): Brain and Language Associations in Autism Spectrum
30 Disorder: A Scoping Review. *J Autism Dev Disord* 52: 725–737.
- 31 28. Morrel J, Singapuri K, Landa RJ, Reetzke R (2023): Neural correlates and
32 predictors of speech and language development in infants at elevated
33 likelihood for autism: a systematic review.
34 <https://doi.org/10.3389/fnhum.2023.1211676>
- 35 29. Crosse MJ, Foxe JJ, Tarrit K, Freedman EG, Molholm S (2022): Resolution of

- 1 impaired multisensory processing in autism and the cost of switching
2 sensory modality [no. 1]. *Commun Biol* 5: 1–17.
- 3 30. Jao Keehn RJ, Sanchez SS, Stewart CR, Zhao W, Grenesko-Stevens EL,
4 Keehn B, Müller R-A (2017): Impaired downregulation of visual cortex
5 during auditory processing is associated with autism symptomatology in
6 children and adolescents with autism spectrum disorder. *Autism Res Off J Int*
7 *Soc Autism Res* 10: 130–143.
- 8 31. Stevenson RA, Siemann JK, Schneider BC, Eberly HE, Woynaroski TG,
9 Camarata SM, Wallace MT (2014): Multisensory temporal integration in
10 autism spectrum disorders. *J Neurosci* 34: 691–7.
- 11 32. Stevenson RA, Segers M, Ferber S, Barense MD, Wallace MT (2014): The
12 impact of multisensory integration deficits on speech perception in children
13 with autism spectrum disorders. *Front Psychol* 5: 379.
- 14 33. Alm M, Behne DM, Wang Y, Eg R (2009): Audio-visual identification of place
15 of articulation and voicing in white and babble noise(a). *J Acoust Soc Am*
16 126: 377–387.
- 17 34. Bertels J, Niesen M, Destoky F, Coolen T, Vander Ghinst M, Wens V, *et al.*
18 (2023): Neurodevelopmental oscillatory basis of speech processing in noise.
19 *Dev Cogn Neurosci* 59: 101181.
- 20 35. Fleming JT, Maddox RK, Shinn-Cunningham BG (2021): Spatial alignment
21 between faces and voices improves selective attention to audio-visual
22 speech. *J Acoust Soc Am* 150: 3085–3100.
- 23 36. Yuan Y, Lleo Y, Daniel R, White A, Oh Y (2021): The Impact of Temporally
24 Coherent Visual Cues on Speech Perception in Complex Auditory
25 Environments. *Front Neurosci* 15.
26 <https://doi.org/10.3389/fnins.2021.678029>
- 27 37. Guiraud JA, Tomalski P, Kushnerenko E, Ribeiro H, Davies K, Charman T, *et*
28 *al.* (2012): Atypical audiovisual speech integration in infants at risk for
29 autism. *PLoS One* 7: e36428.
- 30 38. Lindborg A, Baart M, Stekelenburg JJ, Vroomen J, Andersen TS (2019):
31 Speech-specific audiovisual integration modulates induced theta-band
32 oscillations. *PLoS One* 14: e0219744.
- 33 39. Munhall KG, Gribble P, Sacco L, Ward M (1996): Temporal constraints on the
34 McGurk effect. *Percept Psychophys* 58: 351–362.
- 35 40. van Wassenhove V (2013): Speech through ears and eyes: interfacing the

- 1 senses with the supramodal brain. *Front Psychol* 4. Retrieved August 28,
2 2023, from <https://www.frontiersin.org/articles/10.3389/fpsyg.2013.00388>
- 3 41. Choi I, Lee JY, Lee SH (2018): Bottom-up and top-down modulation of
4 multisensory integration. *Curr Opin Neurobiol* 52: 115–122.
- 5 42. van Wassenhove V, Grant KW, Poeppel D (2005): Visual speech speeds up the
6 neural processing of auditory speech. *Proc Natl Acad Sci U A* 102: 1181–6.
- 7 43. Kushnerenko E, Teinonen T, Volein A, Csibra G (2008): Electrophysiological
8 evidence of illusory audiovisual speech percept in human infants. *Proc Natl*
9 *Acad Sci* 105: 11442–11445.
- 10 44. Musacchia G, Schroeder CE (2009): Neuronal mechanisms, response dynamics
11 and perceptual functions of multisensory interactions in auditory cortex.
12 *Hear Res* 258: 72–79.
- 13 45. Murray MM, Wallace MT (Eds.) (2011): *The Neural Bases of Multisensory*
14 *Processes*. Boca Raton: CRC Press. <https://doi.org/10.1201/9781439812174>
- 15 46. Baum SH, Stevenson RA, Wallace MT (2015): Behavioral, perceptual, and
16 neural alterations in sensory and multisensory function in autism spectrum
17 disorder. *Prog Neurobiol* 134: 140–60.
- 18 47. Foss-Feig JH, Kwakye LD, Cascio CJ, Burnette CP, Kadivar H, Stone WL,
19 Wallace MT (2010): An extended multisensory temporal binding window in
20 autism spectrum disorders. *Exp Brain Res* 203: 381–9.
- 21 48. Kwakye LD, Foss-Feig JH, Cascio CJ, Stone WL, Wallace MT (2011): Altered
22 auditory and multisensory temporal processing in autism spectrum disorders.
23 *Front Integr Neurosci* 4: 129.
- 24 49. Gao M, Lim S, Chubykin AA (2021): Visual familiarity induced 5 Hz
25 oscillations and improved orientation and direction selectivities in V1. *J*
26 *Neurosci*. <https://doi.org/10.1523/JNEUROSCI.1337-20.2021>
- 27 50. Lakatos P, Shah AS, Knuth KH, Ulbert I, Karmos G, Schroeder CE (2005): An
28 Oscillatory Hierarchy Controlling Neuronal Excitability and Stimulus
29 Processing in the Auditory Cortex. *J Neurophysiol* 94: 1904–1911.
- 30 51. Romei V, Brodbeck V, Michel C, Amedi A, Pascual-Leone A, Thut G (2008):
31 Spontaneous Fluctuations in Posterior α -Band EEG Activity Reflect
32 Variability in Excitability of Human Visual Areas. *Cereb Cortex* 18: 2010–
33 2018.
- 34 52. Leszczynski M, Schroeder CE (2019): The Role of Neuronal Oscillations in
35 Visual Active Sensing. *Front Integr Neurosci* 13.

- 1 <https://doi.org/10.3389/fnint.2019.00032>
- 2 53. Poeppel D, Assaneo MF (2020): Speech rhythms and their neural foundations
3 [no. 6]. *Nat Rev Neurosci* 21: 322–334.
- 4 54. Chandrasekaran C, Trubanova A, Stillitano S, Caplier A, Ghazanfar AA
5 (2009): The Natural Statistics of Audiovisual Speech. *PLOS Comput Biol* 5:
6 e1000436.
- 7 55. Arnal LH, Wyart V, Giraud AL (2011): Transitions in neural oscillations
8 reflect prediction errors generated in audiovisual speech. *Nat Neurosci* 14:
9 797–801.
- 10 56. Young AW, Frühholz S, Schweinberger SR (2020): Face and Voice
11 Perception: Understanding Commonalities and Differences. *Trends Cogn Sci*
12 24: 398–410.
- 13 57. Aller M, Økland HS, MacGregor LJ, Blank H, Davis MH (2022): Differential
14 Auditory and Visual Phase-Locking Are Observed during Audio-Visual
15 Benefit and Silent Lip-Reading for Speech Perception. *J Neurosci* 42: 6108–
16 6120.
- 17 58. Hagan CC, Woods W, Johnson S, Calder AJ, Green GGR, Young AW (2009):
18 MEG demonstrates a supra-additive response to facial and vocal emotion in
19 the right superior temporal sulcus. *Proc Natl Acad Sci* 106: 20010–20015.
- 20 59. Hagan CC, Woods W, Johnson S, Green GGR, Young AW (2013):
21 Involvement of Right STS in Audio-Visual Integration for Affective Speech
22 Demonstrated Using MEG. *PLOS ONE* 8: e70648.
- 23 60. Plöchl M, Fiebelkorn I, Kastner S, Obleser J (2022): Attentional sampling of
24 visual and auditory objects is captured by theta-modulated neural activity.
25 *Eur J Neurosci* 55: 3067–3082.
- 26 61. van Wassenhove V, Grant KW, Poeppel D (2007): Temporal window of
27 integration in auditory-visual speech perception. *Neuropsychologia* 45: 598–
28 607.
- 29 62. Stevenson RA, Wallace MT (2013): Multisensory temporal integration: task
30 and stimulus dependencies. *Exp Brain Res* 227: 249–61.
- 31 63. Massaro DW, Cohen MM (1993): The paradigm and the fuzzy logical model of
32 perception are alive and well. *J Exp Psychol Gen* 122: 115–124.
- 33 64. Guillemot P, Graef C, Butters E, Reichenbach T (2023): Audiotactile
34 Stimulation Can Improve Syllable Discrimination through Multisensory
35 Integration in the Theta Frequency Band. *J Cogn Neurosci* 35: 1760–1772.

- 1 65. Power AJ, Mead N, Barnes L, Goswami U (2012): Neural entrainment to
2 rhythmically presented auditory, visual, and audio-visual speech in children.
3 *Front Psychol* 3: 216.
- 4 66. Schroeder CE, Lakatos P, Kajikawa Y, Partan S, Puce A (2008): Neuronal
5 oscillations and visual amplification of speech. *Trends Cogn Sci* 12: 106–13.
- 6 67. Keshavarzi M, Mandke K, Macfarlane A, Parvez L, Gabrielczyk F, Wilson A,
7 Goswami U (2022): Atypical delta-band phase consistency and atypical
8 preferred phase in children with dyslexia during neural entrainment to
9 rhythmic audio-visual speech. *NeuroImage Clin* 35: 103054.
- 10 68. Franchini M, Wood de Wilde H, Glaser B, Gentaz E, Eliez S, Schaer M (2016):
11 Brief Report: A Preference for Biological Motion Predicts a Reduction in
12 Symptom Severity 1 Year Later in Preschoolers with Autism Spectrum
13 Disorders. *Front Psychiatry* 7: 143.
- 14 69. Franchini M, Zoller D, Gentaz E, Glaser B, Wood de Wilde H, Kojovic N, *et*
15 *al.* (2018): Early Adaptive Functioning Trajectories in Preschoolers With
16 Autism Spectrum Disorders. *J Pediatr Psychol* 43: 800–813.
- 17 70. Lord C, Risi S, Lambrecht L, Cook EH, Leventhal BL, DiLavore PC, *et al.*
18 (2000): The Autism Diagnostic Observation Schedule—Generic: A standard
19 measure of social and communication deficits associated with the spectrum
20 of autism. *J Autism Dev Disord* 30: 205–223.
- 21 71. Lord C, Rutter M, DiLavore P, Risi S, Gotham K, Bishop S (2012): Autism
22 diagnostic observation schedule: ADOS-2. *West Psychol. J Psychoeduc*
23 *Assess* 32: 88–92.
- 24 72. *Trotro l'anniversaire de nana* (2013): Storimages.
- 25 73. *Trotro part en vacance* (2013): Storimages.
- 26 74. *Trotro et la boîte a secrets* (2013): Storimages.
- 27 75. *Trotro es tres amoureux* (2013): Storimages.
- 28 76. Olsen A (2012): The Tobii I-VT fixation filter. *Tobii Technol* 21: 4–19.
- 29 77. Boashash B (2015): *Time-Frequency Signal Analysis and Processing: A*
30 *Comprehensive Reference*. Academic Press.
- 31 78. Pelli DG, Tillman KA (2008): The uncrowded window of object recognition.
32 *Nat Neurosci* 11: 1129–35.
- 33 79. Nuthmann A (2013): On the visual span during object search in real-world
34 scenes. *Vis Cogn* 21: 803–837.
- 35 80. Jessen S, Fiedler L, Münte TF, Obleser J (2019): Quantifying the individual

- 1 auditory and visual brain response in 7- month-old infants watching a brief
2 cartoon movie. *NeuroImage*.
3 <https://doi.org/10.1016/j.neuroimage.2019.116060>
- 4 81. Weineck K, Wen OX, Henry MJ (2022): Neural synchronization is strongest to
5 the spectral flux of slow music and depends on familiarity and beat salience
6 ((O. Jensen, B. G. Shinn-Cunningham, & B. Zoefel, editors)). *eLife* 11:
7 e75515.
- 8 82. Cover TM, Thomas JA (2005): Entropy, Relative Entropy, and Mutual
9 Information. *Elements of Information Theory*. John Wiley & Sons, Ltd, pp
10 13–55.
- 11 83. Timme NM, Lapish C (2018): A Tutorial for Information Theory in
12 Neuroscience. *eNeuro* 5. <https://doi.org/10.1523/ENEURO.0052-18.2018>
- 13 84. Delorme A, Makeig S (2004): EEGLAB: an open source toolbox for analysis
14 of single-trial EEG dynamics including independent component analysis. *J*
15 *Neurosci Methods* 134: 9–21.
- 16 85. Crosse MJ, Di Liberto GM, Bednar A, Lalor EC (2016): The Multivariate
17 Temporal Response Function (mTRF) Toolbox: A MATLAB Toolbox for
18 Relating Neural Signals to Continuous Stimuli. *Front Hum Neurosci* 10:
19 604.
- 20 86. Maris E, Oostenveld R (2007): Nonparametric statistical testing of EEG- and
21 MEG-data. *J Neurosci Methods* 164: 177–90.
- 22 87. Dinno A (2015): Nonparametric Pairwise Multiple Comparisons in
23 Independent Groups using Dunn’s Test. *Stata J* 15: 292–300.
- 24 88. Nozaradan S, Peretz I, Missal M, Mouraux A (2011): Tagging the neuronal
25 entrainment to beat and meter. *J Neurosci* 31: 10234–40.
- 26 89. Clouter A, Shapiro KL, Hanslmayr S (2017): Theta Phase Synchronization Is
27 the Glue that Binds Human Associative Memory. *Curr Biol* 27: 3143-
28 3148.e6.
- 29 90. Pefkou M, Arnal LH, Fontolan L, Giraud AL (2017): theta-Band and beta-
30 Band Neural Activity Reflects Independent Syllable Tracking and
31 Comprehension of Time-Compressed Speech. *J Neurosci* 37: 7930–7938.
- 32 91. Berens P (2009): **CircStat** : A *MATLAB* Toolbox for Circular Statistics. *J Stat*
33 *Softw* 31. <https://doi.org/10.18637/jss.v031.i10>
- 34 92. Jochaut D, Lehongre K, Saitovitch A, Devauchelle AD, Olasagasti I, Chabane
35 N, *et al.* (2015): Atypical coordination of cortical oscillations in response to

- 1 speech in autism. *Front Hum Neurosci* 9: 171.
- 2 93. Puschmann S, Daeglau M, Stropahl M, Mirkovic B, Rosemann S, Thiel CM,
3 Debener S (2019): Hearing-impaired listeners show increased audiovisual
4 benefit when listening to speech in noise. *NeuroImage* 196: 261–268.
- 5 94. Haider CL, Suess N, Hauswald A, Park H, Weisz N (2022): Masking of the
6 mouth area impairs reconstruction of acoustic speech features and higher-
7 level segmentational features in the presence of a distractor speaker.
8 *NeuroImage* 252: 119044.
- 9 95. O’CONNOR N, Hermelin B (1965): Sensory dominance: In autistic imbecile
10 children and controls. *Arch Gen Psychiatry* 12: 99–103.
- 11 96. Arnal LH, Morillon B, Kell CA, Giraud A-L (2009): Dual Neural Routing of
12 Visual Facilitation in Speech Processing. *J Neurosci* 29: 13445–13453.
- 13 97. Liu W, Li M, Yi L (2016): Identifying children with autism spectrum disorder
14 based on their face processing abnormality: A machine learning framework.
15 *Autism Res* 9: 888–98.
- 16 98. Sinnett S, Soto-Faraco S, Spence C (2008): The co-occurrence of multisensory
17 competition and facilitation. *Acta Psychol (Amst)* 128: 153–161.
- 18 99. Cuppini C, Ursino M, Magosso E, Rowland BA, Stein BE (2010): An
19 emergent model of multisensory integration in superior colliculus neurons.
20 *Front Integr Neurosci* 4. <https://doi.org/10.3389/fnint.2010.00006>
- 21 100. Irwin J, Harwood V, Kleinman D, Baron A, Avery T, Turcios J, Landi N
22 (2023): Neural and Behavioral Differences in Speech Perception for
23 Children With Autism Spectrum Disorders Within an Audiovisual Context.
24 *J Speech Lang Hear Res* 66: 2390–2403.
- 25 101. Chawarska K, Lewkowicz D, Feiner H, Macari S, Vernetti A (2022):
26 Attention to audiovisual speech does not facilitate language acquisition in
27 infants with familial history of autism. *J Child Psychol Psychiatry* 63: 1466–
28 1476.
- 29 102. Noesselt T, Rieger JW, Schoenfeld MA, Kanowski M, Hinrichs H, Heinze H-
30 J, Driver J (2007): Audiovisual Temporal Correspondence Modulates
31 Human Multisensory Superior Temporal Sulcus Plus Primary Sensory
32 Cortices. *J Neurosci* 27: 11431–11441.
- 33 103. Stacey PC, Kitterick PT, Morris SD, Sumner CJ (2016): The contribution of
34 visual information to the perception of speech in noise with and without
35 informative temporal fine structure. *Hear Res* 336: 17–28.

- 1 104. Han S, Chen Y-C, Maurer D, Shore DI, Lewis TL, Stanley BM, Alais D
2 (2022): The development of audio–visual temporal precision precedes its
3 rapid recalibration. *Sci Rep* 12: 21591.
- 4 105. Venezia JH, Thurman SM, Matchin W, George SE, Hickok G (2016): Timing
5 in audiovisual speech perception: A mini review and new psychophysical
6 data. *Atten Percept Psychophys* 78: 583–601.
- 7 106. Zoefel B, Archer-Boyd A, Davis MH (2018): Phase Entrainment of Brain
8 Oscillations Causally Modulates Neural Responses to Intelligible Speech.
9 *Curr Biol* 28: 401-408 e5.
- 10 107. Henry MJ, Obleser J (2012): Frequency modulation entrains slow neural
11 oscillations and optimizes human listening behavior. *Proc Natl Acad Sci*
12 109: 20095–20100.
- 13 108. Lisman JE, Jensen O (2013): The theta-gamma neural code. *Neuron* 77:
14 1002–16.
- 15 109. Giraud AL, Poeppel D (2012): Cortical oscillations and speech processing:
16 emerging computational principles and operations. *Nat Neurosci* 15: 511–7.
- 17 110. Stefanics G, Hangya B, Hernádi I, Winkler I, Lakatos P, Ulbert I (2010):
18 Phase Entrainment of Human Delta Oscillations Can Mediate the Effects of
19 Expectation on Reaction Speed. *J Neurosci* 30: 13578–13585.
- 20 111. Herbst SK, Stefanics G, Obleser J (2022): Endogenous modulation of delta
21 phase by expectation–A replication of Stefanics et al., 2010. *Cortex* 149:
22 226–245.
- 23
- 24

1 Tables

2

Table 1 The statistical difference among joint model and single models for speech envelope and visual motion

	Dunn's multiple comparisons test	Mean rank diff.	Significant?	Adjusted P Value
ASD	speech envelope	-60.48	Yes	0.0304
	visual motion	-137.9	Yes	<0.0001
TD	speech envelope	-51.09	No	0.1414
	visual motion	-58.47	Yes	0.0373

3

4

Table 2 The statistical difference of the decoding accuracy evolution for speech envelope and visual motion

	Dunn's multiple comparisons test	Mean rank diff.	Significant?	Adjusted P Value
Visual motion	ASD v.s.TD	-37.55	Yes	0.0003
Speech envelope	ASD v.s.TD	-13.32	No	0.8866
TD	V v.s. A	-6.094	No	>0.9999
ASD	V v.s. A	-30.32	Yes	0.0065

A: speech envelope

V: visual motion

5

6

7

Table 3 The statistical difference across groups and frequency bands in stimulus-response coherence

	Dunn's multiple comparisons test	Mean rank diff.	Significant?	Adjusted P Value
ASD A	delta v.s. theta	162.7	Yes	<0.0001
ASD V	delta v.s. theta	191.5	Yes	<0.0001
ASD delta	A v.s. V	-7.774	No	>0.9999
ASD theta	A v.s. V	20.94	No	>0.9999
TD A	delta v.s. theta	178.6	Yes	<0.0001
TD V	delta v.s. theta	176.4	Yes	<0.0001
TD delta	A v.s. V	-0.6563	No	>0.9999
TD theta	A v.s. V	-2.844	No	>0.9999
A delta	ASD v.s. TD	12.3	No	>0.9999
A theta	ASD v.s. TD	28.19	No	>0.9999
V delta	ASD v.s. TD	19.42	No	>0.9999
V theta	ASD v.s. TD	4.406	No	>0.9999

A: speech envelope

V: visual motion

8