

UCSF

UC San Francisco Previously Published Works

Title

Latent neural dynamics encode temporal context in speech.

Permalink

<https://escholarship.org/uc/item/3744x9fq>

Authors

Stephen, Emily

Li, Yuanning

Metzger, Sean

et al.

Publication Date

2023-09-15

DOI

10.1016/j.heares.2023.108838

Peer reviewed



Published in final edited form as:

Hear Res. 2023 September 15; 437: 108838. doi:10.1016/j.heares.2023.108838.

Latent neural dynamics encode temporal context in speech

Emily P Stephen^{a,b}, Yuanning Li^{a,c}, Sean Metzger^a, Yulia Oganian^{a,d}, Edward F Chang^{a,*}

^aDepartment of Neurological Surgery, University of California San Francisco, San Francisco, CA 94143, United States

^bDepartment of Mathematics and Statistics, Boston University, Boston, MA 02215, United States

^cSchool of Biomedical Engineering, ShanghaiTech University, Shanghai, China

^dCenter for Integrative Neuroscience, University of Tübingen, Tübingen, Germany

Abstract

Direct neural recordings from human auditory cortex have demonstrated encoding for acoustic-phonetic features of consonants and vowels. Neural responses also encode distinct acoustic amplitude cues related to timing, such as those that occur at the onset of a sentence after a silent period or the onset of the vowel in each syllable. Here, we used a group reduced rank regression model to show that distributed cortical responses support a low-dimensional latent state representation of temporal context in speech. The timing cues each capture more unique variance than all other phonetic features and exhibit rotational or cyclical dynamics in latent space from activity that is widespread over the superior temporal gyrus. We propose that these spatially distributed timing signals could serve to provide temporal context for, and possibly bind across time, the concurrent processing of individual phonetic features, to compose higher-order phonological (e.g. word-level) representations.

Keywords

Electrocorticography; Superior temporal gyrus; Auditory; Reduced-rank regression; Latent state

This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>)

*Corresponding author., Edward.Chang@ucsf.edu (E.F. Chang).

CRedit authorship contribution statement

Emily P Stephen: Conceptualization, Formal analysis, Software, Visualization, Writing – original draft, Writing – review & editing.

Yuanning Li: Conceptualization, Formal analysis, Visualization, Writing – review & editing. **Sean Metzger:** Conceptualization,

Visualization, Writing – review & editing. **Yulia Oganian:** Conceptualization, Data curation, Formal analysis, Software,

Visualization, Writing – review & editing. **Edward F Chang:** Conceptualization, Funding acquisition, Supervision, Resources,

Writing – review & editing.

Code availability

Custom Python code to perform the iRRR fits is available online (https://github.com/emilyps14/iRRR_python), which is a part of the

Matlab implementation by the original authors (<https://github.com/reagan0323/iRRR>, (Li et al., 2019)). Python code for the analysis

pipeline described above is also available (https://github.com/emilyps14/mtrf_python). We thank Antin and colleagues (Antin et al.,

2021) for their implementation of jPCA in the Python programming language (<https://github.com/bantin/jPCA>), which we used to

perform the jPCA.

Declaration of Competing Interest

The authors declare no competing interests.

Supplementary materials

Supplementary material associated with this article can be found, in the online version, at doi:10.1016/j.heares.2023.108838.

1. Introduction

Natural speech is a continuous stream of complex acoustic features, and listeners build representations of auditory objects at multiple levels, from phonemes, to syllables, words, and phrases (Berwick et al., 2013; Chomsky, 1985). The cortical basis of these dynamic compositional operations is an active area of research. There is evidence that the superior temporal gyrus (STG) performs speech-specific extraction of acoustic-phonetic features (Mesgarani et al., 2014), but where and how these segmental features are composed into longer units like words is less understood. Since the cascade of neural activity evoked by a given acoustic-phonetic feature can last longer than the feature itself (Gwilliams et al., 2020; Khalighinejad et al., 2017; Mesgarani et al., 2014; Näätänen and Picton, 1987; Norman-Haignere et al., 2020), there is potential for overlap in the neural representations over time. Hence the neural computations underlying speech comprehension should have a way to keep track of the temporal context of the individual phonetic units in order to compose them into a higher order unit such as a word (Fischer-Baum, 2018; Gwilliams et al., 2020).

We hypothesized that the mechanisms underlying temporal context tracking and composition in auditory cortex would be reflected in low-dimensional latent dynamics of electrocorticography (ECoG)-scale neural recordings. As neural recordings have grown in dimension, latent state models have become more popular as the explanatory framework for understanding neural computation. We use the terms “latent state” and “latent dynamics” to refer to low-dimensional approximations of high-dimensional neural recordings across time (e.g. recordings across many neurons or many electrodes). For example, principal component analysis (PCA) can be used to reduce a 256-dimensional timeseries of ECoG recordings into a 3-dimensional timeseries (the top 3 principal components) that capture as much variance as possible. PCA is one of many techniques that can be used to capture a high-dimensional signal in low-dimensional terms. In general, if a low-dimensional representation captures important properties of the high-dimensional signal, those properties can often be better described and visualized in low dimensions, for example by plotting a 3-dimensional timeseries as a trajectory in a 3-dimensional plot.

Going further than just plotting latent dynamics, there is a growing trend to use the geometric characteristics of latent states (i.e. the shapes formed by the low-dimensional trajectories) to gain insight into the computational roles that are being played by the network (Russo et al., 2020, 2018; Seely et al., 2016; Vyas et al., 2020). One such geometrical motif is rotational dynamics (Churchland et al., 2012), which happen when the latent dynamics form circles or closed loops. Rotational dynamics may play a computational role in coordinating movements over time in the motor system (Buonomano and Laje, 2010; Cannon and Patel, 2021; Russo et al., 2020, 2018) (see Section 5). While the neural activity underlying speech perception is likely to be very different from that underlying motor sequencing, low-dimensional dynamics across the speech-responsive network in STG could reflect similar computational strategies to coordinate temporal context during speech perception.

There is already reason to believe that STG encodes information about timing: some STG populations respond to amplitude onset events found at the beginning of a sentence after a silent period, or the acoustic edges that occur at the onset of vowels in syllables (called ‘peak rate’ events because they are defined by peaks in the first derivative of the speech envelope timeseries) (Hamilton et al., 2018; Oganian and Chang, 2019). If these signals are strong (representing a large proportion of the variance), temporally similar across different populations, and spatially widespread, they could constitute a meaningful low-dimensional latent state. In fact, Hamilton and colleagues (Hamilton et al., 2018) were able to find low-dimensional dynamics tied to sentence onsets using unsupervised linear dimensionality reduction. Unfortunately, due to the complex nature of the task (with a high-dimensional stimulus space and relevant stimulus features occurring closely in time), unsupervised methods have trouble uncovering dynamics related to other stimulus features, whose neural responses may overlap temporally and spatially with sentence onset responses. This makes it difficult to describe latent dynamics related to peak rate events, which are more closely aligned in timescale to the low-level compositional operations that we seek to describe. Supervised models, on the other hand, have historically focused on individual electrodes and as a result fail to describe latent dynamics that may reflect computational principles on a larger spatial scale.

Here we use a multivariate supervised approach to model the activity across all speech-responsive STG electrodes. Using integrative reduced rank regression (iRRR) (Li et al., 2019), we estimate latent states by reducing the high-dimensional ECoG timeseries into a set of low-dimensional responses to specific stimulus features. In other words, we simultaneously estimate a separate low-dimensional latent state for each stimulus feature, including sentence onsets, peak rate events, and acoustic-phonetic features based on the place and manner of articulation. We find that iRRR outperforms models that treat each electrode individually, indicating that substantial feature-related information is shared across electrodes. The sentence onset and peak rate features explain more of the variance than phonetic features, reaffirming the importance of these timing-related features for encoding in STG. Furthermore, the latent states for the onset and peak rate features are low-dimensional (5 and 6 dimensional, respectively) and distributed over centimeters of cortex, indicating a widespread signal that would be available to coordinate local and downstream processing. Geometrically, the latent dynamics contain a large proportion of rotational dynamics. Projections of the neural responses onto these low-dimensional spaces can be used to decode the time relative to the most recent sentence onset or peak rate event, with performance that is better than decoding from the full high-dimensional responses across all electrodes. We propose that the sentence onset response is an initialization signal and the peak rate latent states encode the time relative to acoustic events at the sentence and syllable scales. For peak rate, this spatially distributed timing signal could be used in local and downstream processing when composing word-level representations from low-level acoustic features.

2. Theory

High gamma amplitudes in neural voltage recordings are known to correlate with the firing rates (Dubey and Ray, 2020; Manning et al., 2009; Ray et al., 2008; Ray and Maunsell, 2011; Scheffer-Teixeira et al., 2013) and dendritic processes (Bédard et al., 2006;

Leszczyński et al., 2020; Miller et al., 2009; Suzuki and Larkum, 2017) of neurons near the electrode (Buzsáki et al., 2012), and we use them here as a proxy for the level of population activity under the ECoG electrodes. Successful previous models of high gamma activity over STG have taken two different approaches: using univariate supervised regression to model single-electrode responses as a function of spectral or linguistic characteristics in the audio speech signal (Aertsen and Johannesma, 1981; Holdgraf et al., 2017; Mesgarani et al., 2014; Oganian and Chang, 2019; Theunissen et al., 2001), and using unsupervised dimensionality reduction to infer latent states from the multivariate signals without reference to the characteristics of the audio stimulus (Hamilton et al., 2018).

2.1. Classic univariate regression modeling

The advantage of regression models is that they characterize the relationship between the neural responses and acoustic features in the speech signal. In classic univariate models, the high gamma responses on individual electrodes are considered to be the result of a convolution of time-dependent receptive fields with corresponding time series of acoustic features. The classic spectrotemporal receptive field (STRF) model (see Section 3.5), for example, uses a mel spectrogram of the stimulus as the acoustic feature representation, resulting in a framework where the neural receptive fields act as a linear filter on the speech spectrogram (Theunissen et al., 2001). Based on the observation that electrode activity over STG reflects information at the level of phonetic features rather than individual phonemes (Mesgarani et al., 2014), Oganian and Chang (2019) used an event-based feature representation to capture these effects and to show that some electrodes additionally have responses triggered by sentence onsets and sharp transients in the acoustic envelope of the speech signal, called peak rate events. While these models have been instrumental in describing the response patterns on individual electrodes, they fail to capture latent dynamics that are shared across multiple electrodes, which could uncover computational principles at work at a larger spatial scale.

2.2. Unsupervised dimensionality reduction modeling

An alternative approach uses unsupervised dimensionality reduction to investigate latent structure in neural responses to speech (Hamilton et al., 2018). Using convex nonnegative matrix factorization, Hamilton and colleagues showed that electrodes can be naturally classified into two groups, “onset” electrodes that have a short increase in high gamma activity at the onset of a sentence, and “sustained” electrodes that show increased high gamma activity throughout the stimulus. This observation is also apparent using principal component analysis (Section 3.12 and Supplementary Figure S1), in which the first component has a characteristic sustained profile, and the second component has the onset profile. Note that the high gamma signals are not intrinsically low-dimensional: 2 dimensions capture only 24% of the variance in speech responsive electrodes (comparable to 16.9% of the variance in all electrodes captured in the first two clusters of (Hamilton et al., 2018)) and 189 dimensions are necessary to capture 80% of the variance. This could be related to the high-dimensional nature of the task: in an unsupervised framework in which the system responds to stimulus features, the response dimensionality needs to be at least as high-dimensional as the task itself (Gao et al., 2017; Stringer et al., 2019). Furthermore, both of these components are time-locked to sentence onset, and it is difficult to connect them or

higher components to other speech features, possibly because the dynamics related to other features are not orthogonal to the sentence-onset subspace or to each other. In particular, the dependence of the neural responses on the peak rate events is not apparent from this analysis, and a model that could capture latent dynamics related to peak rate would be valuable for describing population encoding of shorter timescales.

2.3. Our approach: integrative reduced rank regression

Here we apply a model that combines the advantages of the regression and dimensionality reduction approaches, using multivariate integrative reduced rank regression (iRRR) (Li et al., 2019) to estimate the latent dynamics attributed to each speech feature separately (sentence onsets, peak rate events, and phonetic feature events coded by their place and manner of articulation). This group-reduced-rank model partitions the expected neural activity into a separate latent state for each feature, choosing the best latent dimensionality for each feature while penalizing the total dimensionality across all features. The resulting estimates of feature-specific latent states have explanatory power that goes beyond both individual electrode models and unsupervised dimensionality reduction models.

3. Methods

3.1. Participants

Participants included 11 patients (6M/5F; aged 16–60 years old, median 29) undergoing treatment for intractable epilepsy. As a part of their clinical evaluation for epilepsy surgery, high-density intracranial electrode grids (AdTech 256 channels, 4 mm center-to-center spacing and 1.17 mm diameter) were implanted subdurally over the left peri-Sylvian cortex. All subjects were left-language-dominant (see Table S1 for more clinical and demographic details). All procedures were approved by the University of California, San Francisco Institutional Review Board, and all patients provided informed written consent to participate. Data used in this study was previously reported in (Hamilton et al., 2018).

3.2. Experimental stimuli

Stimuli consisted of 499 English sentences from the TIMIT acoustic-phonetic corpus (Garofolo et al., 1993), spoken by male and female speakers with a variety of North American accents. Stimuli were presented through free-field Logitech speakers at comfortable ambient loudness (~70 dB), controlled by a custom MATLAB script. Participants passively listened to the sentences in 4 blocks, each lasting about 4 min. A subset of 438 sentences were selected for analysis that were heard once by all 11 subjects. The sentences had durations between 0.9 and 2.6 s, with a 400 ms intertrial interval.

3.3. Neural recordings and electrode localization

Neural recordings were acquired at a sampling rate of 3051.8 Hz using a 256-channel PZ2 amplifier or 512-channel PZ5 amplifier connected to an RZ2 digital acquisition system (Tucker-Davis Technologies, Alachua, FL, USA).

Electrodes were localized by coregistering a preoperative T1 MRI scan of the individual subject's brain with a postoperative CT scan of the electrodes in place. Freesurfer was used

to create a 3d model of the individual subjects' pial surfaces, run automatic parcellation to get individual anatomical labels, and warp the individual subject surfaces into the cvs_avg35_inMNI152 average template (Desikan et al., 2006; Fischl et al., 2004). More detailed procedures are described in (Hamilton et al., 2017).

3.4. Preprocessing

For each electrode, the high gamma amplitude time series were extracted from the broadband neural recordings as follows (Edwards et al., 2009; Hamilton et al., 2018; Moses et al., 2016; Oganian and Chang, 2019). First, the signals were downsampled to 400 Hz, rereferenced to the common average in blocks of 16 channels (blocks shared the same connector to the preamplifier), and notch filtered at 60, 120, and 180 Hz to remove line noise and its harmonics. These LFP signals were then filtered using a bank of 8 Gaussian filters with center frequencies logarithmically spaced between 70 and 150 Hz (see Table S2). Using the Hilbert transform, the amplitude of the analytic signal was computed for each of these frequency bands, and for each electrode the high gamma amplitude was defined as the first principal component across these 8 frequency bands. Finally, the high gamma amplitude was further downsampled to 100 Hz and z-scored based on the mean and standard deviation across each experimental block.

3.5. Electrode selection

In order to select speech-responsive electrodes over STG, electrodes were included (1) if they were located over the STG, as identified in the Freesurfer anatomical parcellation of the individual subject cortical surface, and (2) if their high gamma activity was predicted by a linear spectrotemporal model with r^2 above 5% (Hamilton et al., 2018). Note that several electrodes appear to be located away from STG in the cvs_avg35_inMNI152 average template (e.g. Fig. 1a) – this is an artifact of the warping to the average brain.

For this single electrode analysis, the model had the form of a spectrotemporal receptive field (STRF):

$$y(t) = \sum_f \sum_{\tau} s(f, t - \tau) \beta(\tau, f) + e(t) \quad (1)$$

where y is the high gamma amplitude on a single electrode across time t , S is the mel spectrogram of the speech audio signal, β are unknown regression coefficients, and e is the zero-mean Gaussian error term. The frequencies f take on values between 75 Hz and 8 kHz, and delays τ take on values between 0 and 500 ms. By fitting regression coefficients across frequencies and delays, the response on the electrode at a given time is modeled as a function of the recent history of the stimulus spectrogram (up to 500 ms in the past). Ridge regression was used to fit the models (see Section 3.7 for details of the ridge regression framework): the data were split into 80% training and 20% testing data sets, the training data was used to choose the α parameter according to a 5-fold cross-validation, the full training data was fit using the chosen α parameter, and the r^2 was assessed on the testing data (see Section 3.8 for computation of r^2). Electrodes with $r^2 > 0.05$ were included in subsequent

analyses. The selected electrodes and their corresponding r^2 values are shown in Fig. 1A ($N = 331$).

3.6. Regression model setup

The model uses a multivariate adaptation of the event-based regression framework of Oganian and Chang (2019). In matrix form, the model has the following structure:

$$Y = \sum_{f=1}^F X_f B_f + E \quad (2)$$

Where:

- Y is the $T \times N$ matrix of z-scored high gamma amplitude values across electrodes and timepoints. The time dimension represents a concatenation of all 438 sentence stimuli that were heard by every subject, from 500 ms before sentence onset until 500 ms after sentence offset (132,402 timepoints, later split for cross validation, see Section 3.7). The electrode dimension includes speech-responsive electrodes from all subjects (331 electrodes).
- Each $X_f (T \times D)$ represents the delayed feature events for feature f . The first column contains the feature events across time (1 representing an event occurring, 0 otherwise). For peak rate, events were coded by a real-valued magnitude, see Fig. 1B). Following columns contain the same time series, offset by time-delays between 10 ms and 750 ms (76 delays). There were 12 features: sentence onset, peak rate, dorsal, coronal, labial, high, front, low, back, plosive, fricative, and nasal (described below).
- $E (T \times N)$ is Gaussian noise, assumed to be uncorrelated across electrodes
- $B_f (D \times N)$ are the coefficient matrices, i.e. the multivariate temporal response functions (MTRFs), representing the responses of each electrode to the given feature across electrodes and delays
- T : number of timepoints; N : number of electrodes, D : number of delays, F : number of features.

Electrodes from all subjects were included in the same model fit, in keeping with the analysis in (Hamilton et al., 2018), in order to maximize statistical power and spatial coverage of STG. However, the model performance is similar for single subjects (See Section 3.13 and Supplementary Figure S2).

The features used to represent the stimulus were chosen to capture both the phonetic contents of speech, as summarized in (Mesgarani et al., 2014), as well as the speech-envelope landmarks that have been shown to predict neural responses: sentence onsets (Hamilton et al., 2018) and peak rate events (Oganian and Chang, 2019). Sentence onset was defined as the sound onset time for the sentence stimulus. Peak rate was extracted by taking

the derivative of the analytic envelope of the speech signal: the peak rate event times were the times when the derivative reached a maximum, and the peak rate magnitude was the value of the derivative at that time point (Oganian and Chang, 2019). Phonetic feature event times (dorsal, coronal, labial, high, front, low, back, plosive, fricative, nasal) were extracted from time-aligned phonetic transcriptions of the TIMIT corpus, which were timed to the onset of the respective phonemes in the speech signal (Garofolo et al., 1993).

Fig. 1B shows the feature events for an example sentence stimulus, “They’ve never met, you know”. The top two panels show the stimulus waveform and mel spectrogram, respectively, with the times of sentence onset and peak rate events indicated with vertical lines (solid and dashed, respectively). The features fall into two categories: timing (sentence onset and peak rate) and acoustic-phonetic (dorsal, coronal, labial, high, low, front, back, plosive, fricative, nasal). With the exception of peak rate, all of the feature events were encoded as binary time series with a 1 representing an event occurring, and 0 otherwise. For peak rate, the time series contained continuous values representing the slope of the acoustic amplitude signal at the time of maximal change, and 0 at all other times (in Fig. 1B, red lines indicate peak rate event times and red numbers indicate the peak rate magnitude). We chose to include magnitude for peak rate events, because it is known to correlate very well with stressed syllables, i.e. syllables with higher stress will have higher peak rate magnitude.

3.7. Model fitting

We fit the model using ordinary least squares (OLS), ridge regression, and integrative reduced-rank regression (iRRR) (Li et al., 2019). The way we use OLS and ridge regression here is equivalent to traditional univariate modeling, and we include them for comparison to the multivariate iRRR approach. The difference between the three is the objective function that is minimized to choose the fitted coefficient matrices:

$$\{\hat{B}_{f,OLS}\}_{f=1}^F = \underset{B_f \in \mathbb{R}^{D \times N}}{\operatorname{argmin}} \frac{1}{2T} \left\| Y - \sum_{f=1}^F X_f B_f \right\|_{\mathcal{F}}^2 \quad (3)$$

$$\{\hat{B}_{f,ridge}\}_{f=1}^F = \underset{B_f \in \mathbb{R}^{D \times N}}{\operatorname{argmin}} \frac{1}{2T} \left\| Y - \sum_{f=1}^F X_f B_f \right\|_{\mathcal{F}}^2 + \alpha \sum_{f=1}^F \|B_f\|_{\mathcal{F}}^2 \quad (4)$$

$$\{\hat{B}_{f,iRRR}\}_{f=1}^F = \underset{B_f \in \mathbb{R}^{D \times N}}{\operatorname{argmin}} \frac{1}{2T} \left\| Y - \sum_{f=1}^F X_f B_f \right\|_{\mathcal{F}}^2 + \lambda \sum_{f=1}^F w_f \|B_f\|_* \quad (5)$$

where $\|\cdot\|_F$ represents the Frobenius (L2) norm, $\|\cdot\|_*$ represents the nuclear norm (i.e. the sum of the singular values of the bracketed matrix), the w_f s are weights chosen as described below, and α and λ are regularization parameters that are chosen as described below.

The weights used for the iRRR model were chosen to balance the different features (Li et al., 2019):

$$w_f = \sigma(X_f, 1) \{ \sqrt{N} + \sqrt{r(X_f)} \} / T \quad (6)$$

where $\sigma(X_f, 1)$ is the first singular value of the matrix X_f and $r(X_f) = D$ is the rank of matrix X_f . Note that the cost functions Eqs. (3)–(5) treat the noise variance for all electrodes equally – because the high gamma signal on each electrode was z-scored in preprocessing, we assume that the noise variance is the same for all electrodes. In addition, all predictors X_f and responses Y were column-centered before fitting the models.

In iRRR, the nuclear norm penalty acts as an L1 penalty on the singular values of each feature matrix B_f , so the regression tends to find solutions where the feature matrices are low-rank (i.e. sparse in the singular values). Because many of the singular values will be zero, the fitted feature matrices can be represented using a low-dimensional singular value decomposition:

$$\hat{B}_f = U_f S_f V_f^T \quad (7)$$

where U_f is $D \times k$, S_f is $k \times k$, and V_f^T is $k \times N$, for some $k < N$. In other words, the full multivariate feature receptive fields can be represented with a small number of patterns across time (columns of U_f), patterns across electrodes (rows of V_f^T), and corresponding weights (values on the diagonal of S_f). The number of dimensions k can be different for each feature, and it comes from balancing the contribution of the feature to the first term of Eq. (5) (the mean squared error) with the contribution of the feature to the second term (the nuclear norm penalty), relative to other features. Increasing the tuning parameter λ will tend to decrease the total number of dimensions used across all features.

Note that the approach of using a regression framework to fit a group-reduced rank model of neural activity has been used before (Aoi et al., 2020; Aoi and Pillow, 2018): the iRRR framework differs in that it uses an L1 relaxation, resulting in a convex optimization formulation that can be fit efficiently using alternating direction method of multipliers.

In order to compute confidence intervals for model performance metrics (Section 3.8), models were fit using 10-fold cross validation, using group cross validation to keep time points corresponding to the same sentence stimulus in the same fold. For ridge regression and iRRR, an additional nested 5-fold cross validation was used to choose the α and λ parameters within each fold of the outer cross-validation. For ridge regression, a separate

α parameter was chosen for each electrode (consistent with the standard approach for univariate models), while iRRR used a single λ parameter for the full multivariate fit.

3.8. Model performance metrics

Total explained variance (Fig. 1C) was calculated as:

$$r^2 = 1 - \frac{SS_{res}}{SS_{tot}} \quad (8)$$

where the SS_{res} is the residual sum of squares computed on the testing dataset:

$$SS_{res} = \left\| Y - \sum_{f=1}^F X_f B_f \right\|_{\mathcal{F}}^2 \quad (9)$$

and SS_{tot} is the total sum of squares computed on the testing dataset:

$$SS_{tot} = \| Y \|_{\mathcal{F}}^2 \quad (10)$$

The group nuclear norm (Fig. 1D) was computed as the penalty term in the iRRR model:

$$\sum_{f=1}^F w_f \| B_f \|_* \quad (11)$$

Because OLS and ridge regression yield full-rank coefficient matrices, the number of parameters (Fig. 1E) used for both is DN . For iRRR, the number of parameters is $k(D + N + 1)$, based on the singular value decomposition described in Eq. (6).

Unique explained variance for each feature (Fig. 1F) was computed by fitting a reduced iRRR model without the feature f , and then comparing the total explained variance of the full model r_{Full}^2 to the total explained variance of the reduced model r_{-f}^2 . The reduced iRRR model was fit using the same λ value as the full model, chosen using nested cross validation on the full model as described above. For the ‘‘all timing’’ category, the reduced model was fit without sentence onset and peak rate, and for the ‘‘all phonetic’’ category, the reduced model was fit without the phonetic features. The unique explained variance was expressed as a percentage of the full model:

$$100 \times \frac{r_{Full}^2 - r_{-f}^2}{r_{Full}^2} \quad (12)$$

All metrics are reported in terms of the mean across the 10 folds of the cross validation, and 95% confidence intervals are $\pm t_{9,0.975} S / \sqrt{10}$, where s is the sample standard deviation across the 10 cross validation folds. Note that these confidence intervals do not account for the dependence between cross-validation folds due to reuse of samples in training and testing sets, and may therefore be smaller than the true intervals (Austern and Zhou, 2020; Bates et al., 2023; Bengio and Grandvalet, 2004).

Significant differences between conditions were assessed using paired two-tailed t-tests across cross-validation folds (Dietterich, 1998) for the following comparisons (with the resulting p-value ranges):

1. Total explained variance for OLS vs Ridge ($p > 0.05$), OLS vs iRRR ($p < 0.0005$), and Ridge vs iRRR ($p < 0.0005$).
2. Unique explained variance of sentence onset vs each acoustic-phonetic feature and peak rate vs each acoustic-phonetic feature. Here the p-values were Bonferroni corrected across the (2 timing features times 10 acoustic-phonetic features) 20 comparisons. After correction, all comparisons were significant with $p < 0.0005$.
3. Unique explained variance of the combined timing features vs the combined acoustic-phonetic features ($p < 0.0005$).

Similar to the confidence intervals described above, the significance tests did not account for the dependence between cross-validation folds and may therefore have an inflated type II error (Austern and Zhou, 2020; Bates et al., 2023; Bengio and Grandvalet, 2004).

3.9. Computing predicted responses

Given a model fitted with iRRR, the predicted latent response to a stimulus matrix X_f is given by:

$$\hat{Y}_{f,latent} = X_f U_f S_f \quad (13)$$

Where $X_f (T \times D)$ represents the delayed feature events for feature f , U_f is the $D \times k$ time components for feature f , and S_f is a diagonal matrix containing the weights for each component ($k \times k$). $\hat{Y}_{f,latent}$ is a $T \times k$ matrix representing the predicted response within the k -dimensional latent space of the feature. Fig. 3 shows the predicted sentence onset and peak rate responses to the sentence “They’ve never met, you know”.

3.10. jPCA

The plane of fastest rotation for the sentence onset and peak rate latent states (Fig. 3C) was identified by applying jPCA (Churchland et al., 2012) to the feature coefficient matrices \hat{B}_f . Using jPCA, we modeled the temporal receptive fields in the coefficient matrix as a linear dynamical system evolving over delays:

$$\frac{d\hat{B}_f(t)}{dt} = M\hat{B}_f(t) \quad (14)$$

where t indexes the delay dimension of \hat{B}_f , so the dynamical system describes the evolution of an N -dimensional dynamical system over D timepoints. By approximating the derivative on the left hand side using first differences, the transition matrix M can be fit using regression. Furthermore, the purely rotational component of the transition matrix can be isolated by constraining the matrix M to be skew-symmetric, having purely imaginary eigenvalues that come in complex conjugate pairs. The pair of eigenvectors with the largest magnitude eigenvalues describes the plane with the fastest rotations.

It is important to note that jPCA identifies planes with fast rotational dynamics, regardless of whether they capture a large proportion of the variance of the dynamics in the original dynamical system. Classic jPCA uses PCA in preprocessing in order to confine the analysis to six dimensions of largest variance. Here, the iRRR model chooses k dimensions for each feature that are most valuable to the overall fit of the model. Hence there was no need to perform additional PCA to reduce the dimensionality. However, because the coefficient matrices had dimensions capturing very little variance, we did subselect components to capture 98% of the variance of the coefficient matrices. For both sentence onset and peak rate, this corresponded to the top 3 components. Hence the jPCA plane represents the plane of maximal rotation within a 3-dimensional subspace capturing 98% of the variance in the 5-dimensional (or 6-dimensional) coefficient matrix for sentence onset (or peak rate). If we had used more components for the jPCA computation, the rotational dynamics would be stronger but they would capture much less of the variance (using k dimensions vs using 3 dimensions: 2.8% vs 31.8% for sentence onset and 4.8% vs 20.3% peak rate), making them less informative about the overall population dynamics.

Once the jPCs were computed using the coefficient matrices, the predicted trajectory for a given stimulus (Fig. 3F and G) is calculated as:

$$\hat{Y}_{f,jPCA} = X_f J_f \quad (15)$$

$$J_f = [E_1 + E_2, j(E_1 - E_2)]$$

where E_1 and E_2 are the eigenvectors with largest eigenvalues of the skew-symmetric matrix M defined above. J_f is therefore the $N \times 2$ projection matrix from electrode space onto the plane of highest rotation from jPCA.

3.11. Event latency decoding

For the decoding analysis (Fig. 4), a perceptron model was trained to predict the time relative to the most recent feature event (up to 750 ms). The model was designed using the

MLPRegressor class of the sklearn package, with one hidden layer with 20 hidden units using a logistic activation function. We used a simple perceptron model in order to account for possible nonlinearities in the mapping from electrode space / feature latent space to relative times.

Using the same cross-validation framework that was used for iRRR model fitting, the perceptron model was trained using the training data (high gamma amplitudes) either across all electrodes Y or using the projected data onto the latent state subspace:

$$\tilde{Y}_{f,proj} = YV_f \quad (16)$$

where V_f is the $N \times k$ matrix of electrode components for feature f , as above. The $T \times k$ matrix $\tilde{Y}_{f,proj}$ is an approximation of the latent state across time, but it may be contaminated by activity from other features because the V_f matrices do not describe orthogonal subspaces. It also contains activity from noise.

Performance of the models was assessed using r^2 (Eq. (8)) on the held-out testing data for the cross-validation fold. The 95% confidence intervals were computed using the t distribution as described above, and the performance of the models trained on all electrodes was compared to the performance of the models trained on the latent projections using a two-sided paired t-test, as described above (Section 3.8), Bonferroni corrected across the 12 features.

3.12. Principal component analysis

For Supplementary Figure S1, a standard principal component analysis (PCA) was run using the same data matrix as used above (Y , the $T \times N$ matrix of z-scored high gamma amplitude values across 331 electrodes during the presentation of 438 sentences). Because the data are already centered, PCA is just a singular value decomposition of the data matrix:

$$Y = U_Y S_Y V_Y^T \quad (17)$$

where U_Y is a $T \times N$ orthogonal matrix where the columns represent the N principal components across time, and V_Y^T is a $N \times N$ matrix where the rows represent the spatial support of each principal component. S_Y is a diagonal $N \times N$ matrix with ordered diagonal elements s_1, s_2, \dots, s_N . The percent of the variance explained by a component i can be calculated as:

$$100 \times \frac{s_i^2}{\sum_{j=1}^N s_j^2} \quad (18)$$

Note that the “percent explained variance” in PCA (Supplementary Figure S1) is not comparable to the “total explained variance” in the regression analysis (Section 3.8 and Fig. 1C), because the PCA explained variance is computed on the training data, while the regression explained variance is computed on held-out testing data.

3.13. Single subject analysis

For Supplementary Figure S2, the entire pipeline was run for a single subject, SL04. In this analysis, the Y matrix defined above of z-scored high gamma amplitude values across electrodes and timepoints was restricted to only the speech-responsive STG electrodes from SL04 (45 electrodes). This subject was chosen based on the large number of speech-responsive electrodes over STG, and their coverage of both middle and posterior STG. Surface plots in Figure S2 use the subject’s cortical surface, without warping to the average brain.

4. Results

The fits to our integrative reduced rank regression model reveal that high gamma responses to speech stimuli across hundreds of electrodes can be parsimoniously represented as a combination of a few low-dimensional latent state responses to specific feature events in the stimulus. Two latent states in particular, corresponding to the sentence onset and peak rate features, reflect a large proportion of the explained variance in the model, and their dynamic properties suggest specific computational roles in the speech perception network.

4.1. iRRR outperforms models that treat each electrode individually, and sentence onset and peak rate capture more of the variance than phonetic features

Fig. 1C–E compare the three different fitting frameworks: OLS, ridge regression, and iRRR. Because the regression framework is the same for all three, the fitted models have very similar total explained variance (r^2 computed over all electrodes, Fig. 1C). All of the models have a proportion of explained variance of about 0.2, which can be partially explained by the fact that each stimulus was presented only once, so the data contains both stimulus-related activity and trial-specific noise. In addition, this is an aggregate over all speech-responsive electrodes: some electrodes are more stimulus-driven than others (see Fig. 1A). Comparing the three fitting frameworks, iRRR by design achieves a much smaller nuclear norm (Fig. 1D), which results in solutions that can be described with 94% fewer parameters than OLS and ridge regression (Fig. 1E). The fact that the iRRR model captures as much information as the single-electrode models using far fewer parameters suggests that substantial feature-related information is shared across electrodes.

Fig. 1F shows the unique explained variance of each of the features in the iRRR fit: sentence onset and peak rate explain a larger percentage of the full model variance than each of the phonetic features ($p < 0.0005$ for all comparisons using a two-sided paired t -test after Bonferroni correction). This suggests that these two timing features reflect a substantial amount of the speech-induced response across STG.

When the features are grouped into timing (sentence onset and peak rate) and phonetic (all other features) groups, both groups explain a large proportion of the variance (15%

and 22%, respectively). Comparing the groups, however, the phonetic features explain more of the unique variance than the timing features ($p < 0.0005$, two-sided paired t -test). This could be surprising in light of the individual feature comparisons: while timing features capture more explained variance than phonetic features when compared individually, when combined they capture less explained variance. This is likely due to (1) correlations between individual phonetic features that lead to lower individual unique explained variance and (2) the fact that more electrodes respond to sentence onset and peak rate than individual phonetic features (Oganian and Chang, 2019), meaning that sentence onset and peak rate have more widespread spatial support than the more spatially localized phonetic features. This more widespread spatial support means that the iRRR model is better able to consolidate the activity patterns across multiple electrodes, i.e. capture the latent dynamics, for the sentence onset and peak rate features than for the phonetic features. Accordingly, the following two sections describe the latent state representations for the sentence onset and peak rate features in more detail.

4.2. The model fit captures known response differences between pSTG and mSTG

In Hamilton and colleagues' (Hamilton et al., 2018) unsupervised model, the "onset" cluster of electrodes was found to occur primarily over the posterior portion of STG (pSTG). This observation led them to propose that pSTG may play a role in detecting temporal landmarks at the sentence and phrase level, because the short-latency, short-duration responses to sentence onsets in pSTG would be able to encode the event time with high temporal resolution. This idea fits well within a long history of evidence that stimulus responses in mSTG have longer latencies and longer durations than those in pSTG (Hamilton et al., 2021; Jasmin et al., 2019; Yi et al., 2019). Here, the model fits recapitulate these known differences between mSTG and pSTG.

As discussed above (Eq. (6)), the feature response matrices that are fitted by the iRRR model can be decomposed into a small number of components across time ("time components", columns of U_f), components across electrodes ("spatial components", rows of V_f^T), and corresponding weights (values on the diagonal of S_f). Fig. 2 shows the Sentence Onset and Peak Rate fitted feature matrices decomposed in this way (Since U_f and V_f are orthonormal, their columns are unit vectors: as a result, their units are arbitrary and can be best interpreted in relative terms).

Fig. 2A and B show the time components scaled by their corresponding weights, and Fig. 2C and D show the first two spatial components. To illustrate how the low dimensional components map back to the response functions for individual electrodes, Fig. 2E and F show the individual electrode response functions (rows of \hat{B}_f), colored by the spatial component from Fig. 2C and D.

Looking at the left panel of Fig. 2C and 2E, we can see that electrodes that have large values in the first spatial component (red circles in Fig. 2C, left) have relatively larger overall responses to sentence onset events (red lines in Fig. 2E, left). These electrodes occur primarily over pSTG (i.e. posterior to the lateral exit point of the transverse temporal sulcus), which is in line with previous findings (Hamilton et al., 2018).

For peak rate, the first component plays the same role: electrodes that have larger values in the first spatial component (Fig. 2D, left) have relatively larger overall responses to peak rate events (Fig. 2F, left). Electrodes with large peak rate responses are not limited to pSTG like sentence onset electrodes: rather, they are distributed over all of STG. In other words, the encoding of peak rate in STG is not focal but is distributed over centimeters of cortex, suggesting a representation on a large spatial scale. Interestingly, the second component does appear to have a spatial distinction between pSTG and mSTG: electrodes with positive values for the second component tend to occur over pSTG, while electrodes with negative values for the second component tend to occur over mSTG (i.e. anterior to the lateral exit point of the transverse temporal sulcus, Fig. 2D, right). The negative and positive values distinguish response functions by their temporal response profile: positive values correspond to electrodes that have an early peak rate response, while negative values correspond to electrodes that have a late peak rate response (Fig. 2F, right). This suggests that peak rate responses over pSTG are faster than peak rate responses over mSTG.

4.3. Feature latent states have rotational dynamics that capture continuous relative timing information

To show how the latent states behave during the presentation of a stimulus, we used the fitted model to predict the dynamics in each latent state during the presentation of the sentence “They’ve never met, you know” (Fig. 3, see Section 3.9 for the calculation of the predicted responses).

The sentence onset latent space has 5 dimensions and the peak rate latent space has 6 dimensions. While the sentence onset feature only occurs once at the beginning of the stimulus, evoking a single response across the sentence onset dimensions, the peak rate feature occurs several times, and the dynamics of the peak rate latent state do not go back to baseline in between peak rate events (Fig. 3B and C). Plotting the top three dimensions, which capture more than 98% of the variance in the coefficient matrices (\hat{B}_t), shows cyclical dynamics for both sentence onset and peak rate (Fig. 3D and E): the sentence onset state rotates once at the beginning of the sentence, and the peak rate latent state rotates 3–4 times, once after each peak rate event.

To quantify this effect, we used jPCA (Churchland et al., 2012) to identify the most rotational 2 dimensional subspace within the top three components of \hat{B}_t . These planes capture 31.8% and 20.3% of the variance in the sentence onset and peak rate coefficient matrices, respectively, and they highlight the cyclical dynamics that were visible in the top 3 dimensions (Fig. 3F and G).

Note that seeing cyclical dynamics in the latent states is not necessarily surprising: the coefficient matrices \hat{B}_t describe smooth multivariate evoked responses that will tend to start and end at the same baseline. Indeed, the cyclical dynamics may reflect a so-called “horseshoe effect” arising from short- and long-latency responses to the same events (Elsayed and Cunningham, 2017; Michaels et al., 2016), as is evident in Fig. 2F. Our data and model are also not intended to distinguish between a dynamical code versus a representational code, which is an ongoing controversy in the field: a representational code

explains neural activity with behavioral or external factors, while a dynamical code explains neural activity as a function of previous neural activity. While these two frameworks are not mutually exclusive, neural systems may be better explained by one or the other in different situations (Michaels et al., 2016; Russo et al., 2018; Vyas et al., 2020).

Here, we highlight the rotational dynamics to motivate a geometrical argument for the role of the peak rate responses in downstream processing. We will make the case (see Section 5) that the structure of the peak rate responses enables them to act as a temporal context signal against which other features are organized. In order for the peak rate latent state to play this role, the trajectories should be sufficiently spread out in latent space to enable downstream areas to decode the time relative to the most recent peak rate event using just the instantaneous latent state. We investigate whether this is true in the next section.

4.4. Latent states from the model can be used to decode time relative to feature events

So far, we have described how the model is fit using known feature event times, and how the fitted model can be used to predict responses given new feature events. We also wanted to know whether the model fit could be used to decode the timing of events, which would indicate that sufficient information is contained in the feature responses for downstream areas to use them as temporal context signals.

The set of spatial components for each feature defines a feature-specific subspace of the overall electrode space. The projection of the observed high gamma time series onto this subspace is an approximation of the feature latent state (note that it is not exact, because the different feature subspaces are not orthogonal to each other). We asked whether this latent projection time series could be used to decode the time since the most recent feature event.

Fig. 4 shows the result of this analysis (details of the methods are in Section 3.11): a perceptron model was trained to decode the time since the most recent feature event up to 750 ms, given either the activity on the full set of electrodes or the projection of the electrode activity onto the corresponding feature subspace. The decoder for sentence onset performs slightly better when using all electrodes, which may be due to the large proportion of the overall activity that is time-locked to sentence onsets (see Supplementary Figure S1). For all other features, however, decoder performance using the reduced-dimensional latent subspaces performs even better than decoding using the full dimensional activity across electrodes (paired *t*-test over 10 cross validation folds, $p < 0.05$ with Bonferroni correction across 12 features). Because no information is gained in the projection operation, this is an indication that projecting onto the latent subspaces increases the signal to noise ratio, i.e. removes activity that is irrelevant to decoding relative time.

5. Discussion

We have shown that a low dimensional regression model, iRRR, performs as well as classic models in representing high-gamma responses to timing and phonetic features of auditory stimuli, while using far fewer parameters. It accomplishes this compression by capturing similarities in feature responses that are shared across electrodes, which enables a low-dimensional latent state interpretation of the dynamics of high gamma responses to

stimulus features. The sentence onset and peak rate features capture more unique variance than the other (phonetic) features, their responses are spread over both mSTG and pSTG, and their latent states show rotational dynamics that repeat after each event. Based on the geometry, duration, and spatial extent of the latent dynamics, we make the case that the sentence onset response could act as an initialization signal to kick the network into a speech-encoding state, while the peak rate response could provide a widespread temporal context signal that could be used to compose word-level representations from low-level acoustic and phonetic features.

The large magnitude of sentence onset responses in ECoG high gamma responses has been reported before (Hamilton et al., 2018): here, we confirm their large contribution to STG responses both using our iRRR model (Fig. 1) and using PCA (Supplementary Figure S1). Importantly, the latent dynamics related to sentence onset last about 600 ms (Fig. 2a). Since sentences in English often last longer than 600 ms (e.g. the sentences in the TIMIT corpus used here ranged from 900 ms to 2.6 s), these onset-related dynamics are unsuited to encode temporal context on an entire sentence level. Furthermore, sentence boundaries in continuous natural speech are rarely indicated with pauses or silence (Yoon et al., 2007), meaning that neural responses to acoustic onsets are unlikely to code sentence transitions. Rather, the latent dynamics in response to onsets may serve as a non-speech specific temporal indicator of the transition from silence to sound, occurring during perception of any auditory stimulus. During speech perception, the speech-related cortical networks could use this non-specific event as a reset or initialization signal. The idea that a large transient in the latent state could act to transition a network between states is also thought to occur in the motor system, where condition-invariant movement onset responses in the latent state mark the transition from motor preparation to motor behavior (Kaufman et al., 2016).

With regard to the peak rate dynamics, we propose that the computational role of the peak rate feature response is to keep track of word-level temporal context using a clock-like representation. The idea that structured latent state dynamics can act as clocks has been proposed in several different cognitive domains, most commonly in the motor system (Buonomano and Laje, 2010; Churchland et al., 2012; Remington et al., 2018; Vyas et al., 2020) (c.f. (Lebedev et al., 2020)) and in temporal interval estimation and perception (Cannon and Patel, 2021; Gámez et al., 2019; Mauk and Buonomano, 2004; Wang et al., 2018). In the motor system, Russo and colleagues (Russo et al., 2020) describe population dynamics in primary motor cortex (M1) and supplementary motor area (SMA) while a monkey performed a cyclic motor action. The population dynamics in M1 were rotational, exhibiting one rotation for each motor cycle, while the dynamics in SMA were shaped like a spiral, where 2-dimensional rotations for each motor cycle were translated along a third dimension. They proposed that this structure would be well-suited to keep track of progress through multi-cycle actions: each rotation encodes a single action, and translation along the third dimension encodes progress through the motor sequence. The rotational component of SMA population trajectories has also been suggested to operate as a time-keeping signal in auditory beat perception, where rotations through latent space keep track of the interval between beats (Cannon and Patel, 2021).

The peak rate latent state in STG could similarly be playing a computational role in auditory speech perception: the rotations in the peak rate subspace could serve to keep track of the time relative to the peak rate event, chunking time into intervals starting at the onset of a vowel. These intervals could then be used by downstream processing to give temporal context to the fine-grained phonetic feature information conveyed by other subpopulations. In other words, the rotational peak rate latent state could provide a temporal scaffolding on which individual phonetic features can be organized. Fig. 5 illustrates this idea: when hearing the sentence “It had gone like clockwork,” the peak rate latent state partitions the sentence into four rotations, each one capturing the time since the most recent peak rate event. Downstream processing streams could combine this information with the phonetic feature information to put the phonetic feature events into their local context, here at the level of words or small sets of words (Fig. 5C). Peak rate is in a unique position to play this role: it is the only feature that repeats within the linguistic structure of speech at the level of syllables/words, without reference to the linguistic contents. In addition, the peak rate responses are distributed over centimeters of cortex (Fig. 2D) so the temporal context information would be widely available to local and downstream processing.

In order for the peak rate latent state to play this role, it should have a couple of properties. First, there should be a mapping from points in state space to different relative times. As we showed in Fig. 3, the rotational dynamics cause different relative times to be encoded in different locations of the latent space. Second, the trajectories in latent space should be consistent enough to support decoding of relative time in the presence of noise. In Fig. 4, we showed that the projections of the neural activity onto the subspaces spanned by the feature latent states support decoding of the time relative to the most recent feature event. Note that while the latent state projections support decoding better than decoding from the full high-dimensional signal, the actual performance for peak rate is somewhat low (~50%). A possible reason for this could be that some peak rate events are more effective at driving the latent state than others (even after accounting for peak rate magnitude, as the model does), resulting in inconsistent decoding of the time since the most recent peak rate event.

Beyond the two-dimensional rotational dynamics, the peak rate latent trajectory forms a spiral in 3 dimensions (Fig. 5B), similar to population trajectories in SMA during motor sequences (Russo et al., 2020). This suggests that the peak rate subpopulation may additionally encode the ordering of the word-level intervals within a larger linguistic context, such as the phrase level.

Furthermore, the representation of these intervals does not require top-down predictive coding (Hovsepyan et al., 2020; Lewis and Bastiaansen, 2015; Park et al., 2015; Pefkou et al., 2017) or entrainment of ongoing oscillations (Canolty, 2007; Ghitza, 2011; Giraud and Poeppel, 2012; Hovsepyan et al., 2020; Martin, 2020; Pittman-Polletta et al., 2020): in our model they are implemented via event-related potentials triggered by discrete acoustic (peak rate) events. While top-down and oscillatory mechanisms may play important roles in speech perception, our model demonstrates that some speech segmentation and context processing can be performed without them.

The events that we focus on for speech segmentation are peak rate events, moments of sharp increases in the acoustic envelope. The peak rate events in the model are coded with their magnitude (the slope of the rise in the acoustic envelope), which allows the model dynamics to change proportionally to the size of the event. This is important because peak rate events, also called auditory onset edges (Biermann and Heil, 2000; Doelling et al., 2014; Heil and Neubauer, 2001), differ in magnitude based on the stress level of the corresponding syllable (Oganian and Chang, 2019). This means that the dynamics triggered by peak rate events are sensitive to prosodic structure, both stressed syllables within words and stressed words within phrases. To investigate this further, it would be helpful to use a speech stimulus corpus with more complex prosodic structure than the TIMIT corpus used here.

In summary, our model (iRRR) represents STG high gamma responses to natural speech stimuli as a superposition of responses to individual phonetic and timing features, where each feature has a corresponding low-dimensional latent state that is shared across electrodes. It performs as well as single electrode models while using far fewer parameters, indicating that substantial feature-related information is shared across electrodes. Sentence onset and peak rate events, features representing timing at the sentence and syllable scales, capture more unique variance than phonetic features. The latent dynamics for sentence onset and peak rate contain information about the time since the most recent (sentence onset or peak rate) event, and the information is distributed across centimeters of cortex. We make the case that for peak rate, this relative timing information could play a role in composing word-level representations from low-level acoustic features, without requiring oscillatory or top-down mechanisms.

Supplementary Material

Refer to Web version on PubMed Central for supplementary material.

Acknowledgments

This work was supported by grants from the NIH (R01-DC012379 and U01-NS117765 to EFC). This research was also supported by Bill and Susan Oberndorf, The Joan and Sandy Weill Foundation, and The William K. Bowes Foundation. The authors would like to thank the members of the Chang lab at UCSF as well as James Hieronymus and Benjamin Antin for valuable feedback.

Data availability

Data will be made available on request.

References

- Aertsen AMHJ, Johannesma PIM, 1981. The spectro-temporal receptive field: a functional characteristic of auditory neurons. *Biol. Cybern* 42, 133–143. doi:10.1007/BF00336731. [PubMed: 7326288]
- Antin B, Shenoy K, Linderman S, 2021. Probabilistic jPCA: a constrained model of neural dynamics, in: *Cosyne Abstracts 2021*. Presented at the Cosyne21, Online.
- Aoi MC, Mante V, Pillow JW, 2020. Prefrontal cortex exhibits multidimensional dynamic encoding during decision-making. *Nat. Neurosci* 23, 1410–1420. doi:10.1038/s41593-020-0696-5. [PubMed: 33020653]

- Aoi M and Pillow JW, 2018. Model-based targeted dimensionality reduction for neuronal population data. *Advances in neural information processing systems*, 31.
- Austern M and Zhou W, 2020. Asymptotics of cross-validation. arXiv preprint arXiv:2001.11111.
- Bates S, Hastie T, Tibshirani R, 2023. Cross-validation: what does it estimate and how well does it do it? *J. Am. Stat. Assoc* 1–12.
- Bédard C, Kröger H, Destexhe A, 2006. Does the 1/f frequency scaling of brain signals reflect self-organized critical states? *Phys. Rev. Lett* 97, 118102. doi:10.1103/PhysRevLett.97.118102. [PubMed: 17025932]
- Bengio Y, Grandvalet Y, 2004. No unbiased estimator of the variance of k-fold cross-validation. *J. Mach. Learn. Res* 5, 1089–1105.
- Berwick RC, Friederici AD, Chomsky N, Bolhuis JJ, 2013. Evolution, brain, and the nature of language. *Trends Cogn. Sci* 17, 89–98. doi:10.1016/j.tics.2012.12.002. [PubMed: 23313359]
- Biermann S, Heil P, 2000. Parallels between timing of onset responses of single neurons in cat and of evoked magnetic fields in human auditory cortex. *J. Neurophysiol* 84, 2426–2439. doi:10.1152/jn.2000.84.5.2426. [PubMed: 11067985]
- Buonomano DV, Laje R, 2010. Population clocks: motor timing with neural dynamics. *Trends Cogn. Sci* 14, 520–527. doi:10.1016/j.tics.2010.09.002. [PubMed: 20889368]
- Buzsáki G, Anastassiou CA, Koch C, 2012. The origin of extracellular fields and currents — EEG, ECoG, LFP and spikes. *Nat. Rev. Neurosci* 13, 407–420. doi:10.1038/nrn3241. [PubMed: 22595786]
- Cannon JJ, Patel AD, 2021. How beat perception co-opts motor neurophysiology. *Trends Cogn. Sci* 25, 137–150. doi:10.1016/j.tics.2020.11.002. [PubMed: 33353800]
- Canolty RT, 2007. Spatiotemporal dynamics of word processing in the human brain. *Front. Neurosci* 1, 185–196. doi:10.3389/neuro.01.1.1.014.2007. [PubMed: 18982128]
- Chomsky N, 1985. Syntactic structures. *Mouton de Gruyter*.
- Churchland MM, Cunningham JP, Kaufman MT, Foster JD, Nuyujukian P, Ryu SI, Shenoy KV, 2012. Neural population dynamics during reaching. *Nature* 1–8. doi:10.1038/nature11129.
- Desikan RS, Ségonne F, Fischl B, Quinn BT, Dickerson BC, Blacker D, Buckner RL, Dale AM, Maguire RP, Hyman BT, Albert MS, Killiany RJ, 2006. An automated labeling system for subdividing the human cerebral cortex on MRI scans into gyral based regions of interest. *Neuroimage* 31, 968–980. doi:10.1016/j.neuroimage.2006.01.021. [PubMed: 16530430]
- Dietterich TG, 1998. Approximate statistical tests for comparing supervised classification learning algorithms. *Neural Comput* 10, 1895–1923. doi:10.1162/089976698300017197. [PubMed: 9744903]
- Doelling K, Arnal L, Ghizza O, Poeppel D, 2014. Acoustic landmarks drive delta-theta oscillations to enable speech comprehension by facilitating perceptual parsing. *Neuroimage* 85. doi:10.1016/j.neuroimage.2013.06.035.
- Dubey A, Ray S, 2020. Comparison of tuning properties of gamma and high-gamma power in local field potential (LFP) versus electrocorticogram (ECoG) in visual cortex. *Sci. Rep* 10, 5422. doi:10.1038/s41598-020-61961-9. [PubMed: 32214127]
- Edwards E, Soltani M, Kim W, Dalal SS, Nagarajan SS, Berger MS, Knight RT, 2009. Comparison of time–frequency responses and the event-related potential to auditory speech stimuli in human cortex. *J. Neurophysiol* 102, 377–386. doi:10.1152/jn.90954.2008. [PubMed: 19439673]
- Elsayed GF, Cunningham JP, 2017. Structure in neural population recordings: an expected byproduct of simpler phenomena? *Nat. Neurosci* 20, 1310–1318. doi:10.1038/nn.4617. [PubMed: 28783140]
- Fischer-Baum S, 2018. A Common Representation of Serial Position in Language and Memory. Elsevier, pp. 31–54. doi:10.1016/bs.plm.2018.08.002 *Psychology of Learning and Motivation*.
- Fischl B, van der Kouwe A, Destrieux C, Halgren E, Ségonne F, Salat DH, Busa E, Seidman LJ, Goldstein J, Kennedy D, Caviness V, Makris N, Rosen B, Dale AM, 2004. Automatically parcellating the human cerebral cortex. *Cereb. Cortex* 14, 11–22. doi:10.1093/cercor/bhg087. [PubMed: 14654453]
- Gámez J, Mendoza G, Prado L, Betancourt A, Merchant H, 2019. The amplitude in periodic neural state trajectories underlies the tempo of rhythmic tapping. *PLOS Biol.* 17, e3000054. doi:10.1371/journal.pbio.3000054. [PubMed: 30958818]

- Gao P, Trautmann E, Yu B, Santhanam G, Ryu S, Shenoy K, Ganguli S, 2017. A theory of multineuronal dimensionality, dynamics and measurement. *BioRxiv* 214262.
- Garofolo JS, Lamel LF, Fisher WM, Pallett DS, Dahlgren NL, Zue V, Fiscus JG, 1993. TIMIT Acoustic-Phonetic Continuous Speech Corpus LDC93S1. Web Download Philadelphia: Linguistic Data Consortium, 1993. doi:10.35111/17gk-bn40.
- Ghitza O, 2011. Linking speech perception and neurophysiology: speech decoding guided by cascaded oscillators locked to the input rhythm. *Front. Psychol* 2. doi:10.3389/fpsyg.2011.00130.
- Giraud AL, Poeppel D, 2012. Cortical oscillations and speech processing: emerging computational principles and operations. *Nat. Neurosci* 15, 511–517. doi:10.1038/nn.3063. [PubMed: 22426255]
- Gwilliams L, King JR, Marantz A, Poeppel D, 2020. Neural dynamics of phoneme sequencing in real speech jointly encode order and invariant content (preprint). *Neuroscience* doi:10.1101/2020.04.04.025684.
- Hamilton LS, Chang DL, Lee MB, Chang EF, 2017. Semi-automated anatomical labeling and inter-subject warping of high-density intracranial recording electrodes in electrocorticography. *Front. Neuroinform* 11. doi:10.3389/fninf.2017.00062.
- Hamilton LS, Edwards E, Chang EF, 2018. A spatial map of onset and sustained responses to speech in the human superior temporal gyrus. *Curr. Biol* 28. doi:10.1016/j.cub.2018.04.033, 1860–1871.e4. [PubMed: 29861132]
- Hamilton LS, Oganian Y, Hall J, Chang EF, 2021. Parallel and distributed encoding of speech across human auditory cortex. *Cell* 184 (18), 4626–4639. [PubMed: 34411517]
- Heil P, Neubauer H, 2001. Temporal integration of sound pressure determines thresholds of auditory-nerve fibers. *J. Neurosci* 21, 7404–7415. doi:10.1523/JNEUROSCI.21-18-07404.2001. [PubMed: 11549751]
- Holdgraf CR, Rieger JW, Micheli C, Martin S, Knight RT, Theunissen FE, 2017. Encoding and decoding models in cognitive electrophysiology. *Front. Syst. Neurosci* 11. doi:10.3389/fnsys.2017.00061.
- Hovsepian S, Olasagasti I, Giraud AL, 2020. Combining predictive coding and neural oscillations enables online syllable recognition in natural speech. *Nat. Commun* 11, 3117. doi:10.1038/s41467-020-16956-5. [PubMed: 32561726]
- Jasmin K, Lima CF, Scott SK, 2019. Understanding rostral–caudal auditory cortex contributions to auditory perception. *Nat. Rev. Neurosci* 20, 425–434. doi:10.1038/s41583-019-0160-2. [PubMed: 30918365]
- Kaufman MT, Seely JS, Sussillo D, Ryu SI, Shenoy KV, Churchland MM, 2016. The largest response component in the motor cortex reflects movement timing but not movement type. *eNeuro* 3. doi:10.1523/ENEURO.0085-16.2016, ENEURO.0085–16.2016.
- Khalighinejad B, Cruzatto da Silva G, Mesgarani N, 2017. Dynamic encoding of acoustic features in neural responses to continuous speech. *J. Neurosci* 37, 2176–2185. doi:10.1523/JNEUROSCI.2383-16.2017. [PubMed: 28119400]
- Lebedev MA, Ninenko I, Ossadtchi A, 2020. Rotational dynamics versus sequence-like responses (preprint). *Neuroscience* doi:10.1101/2020.09.16.300046.
- Leszczynski M, Barczak A, Kajikawa Y, Ulbert I, Falchier AY, Tal I, Haegens S, Melloni L, Knight RT, Schroeder CE, 2020. Dissociation of broadband high-frequency activity and neuronal firing in the neocortex. *Sci. Adv* 6, eabb0977. doi:10.1126/sciadv.abb0977. [PubMed: 32851172]
- Lewis AG, Bastiaansen M, 2015. A predictive coding framework for rapid neural dynamics during sentence-level language comprehension. *CORTEX* 68, 155–168. doi:10.1016/j.cortex.2015.02.014.
- Li G, Liu X, Chen K, 2019. Integrative multi-view regression: bridging group-sparse and low-rank models. *Biometrics* 75, 593–602. doi:10.1111/biom.13006. [PubMed: 30456759]
- Manning JR, Jacobs J, Fried I, Kahana MJ, 2009. Broadband shifts in local field potential power spectra are correlated with single-neuron spiking in humans. *J. Neurosci* 29, 13613–13620. doi:10.1523/JNEUROSCI.2041-09.2009. [PubMed: 19864573]
- Martin AE, 2020. A compositional neural architecture for language. *J. Cogn. Neurosci* 32, 1407–1427. doi:10.1162/jocn_a_01552. [PubMed: 32108553]

- Mauk MD, Buonomano DV, 2004. The neural basis of temporal processing. *Annu. Rev. Neurosci* 27, 307–340. doi:10.1146/annurev.neuro.27.070203.144247. [PubMed: 15217335]
- Mesgarani N, Cheung C, Johnson K, Chang EF, 2014. Phonetic Feature Encoding in Human Superior Temporal Gyrus. *Science* 343, 1006–1010. doi:10.1126/science.1245994. [PubMed: 24482117]
- Michaels JA, Dann B, Scherberger H, 2016. Neural population dynamics during reaching are better explained by a dynamical system than representational tuning. *PLoS Comput. Biol* 12, e1005175. [PubMed: 27814352]
- Miller KJ, Sorensen LB, Ojemann JG, Den Nijs M, 2009. Power-law scaling in the brain surface electric potential. *PLoS Comput. Biol* 5, e1000609. doi:10.1371/journal.pcbi.1000609.g005. [PubMed: 20019800]
- Moses DA, Mesgarani N, Leonard MK, Chang EF, 2016. Neural speech recognition: continuous phoneme decoding using spatiotemporal representations of human cortical activity. *J. Neural Eng* 13, 056004. doi:10.1088/1741-2560/13/5/056004. [PubMed: 27484713]
- Näätänen R, Picton T, 1987. The N1 wave of the human electric and magnetic response to sound: a review and an analysis of the component structure. *Psychophysiology* 24, 375–425. doi:10.1111/j.1469-8986.1987.tb00311.x. [PubMed: 3615753]
- Norman-Haignere SV, Long LK, Devinsky O, Doyle W, Irobunda I, Merricks EM, Feldstein NA, McKhann GM, Schevon CA, Flinker A, Mesgarani N, 2020. Multiscale integration organizes hierarchical computation in human auditory cortex (preprint). *Neuroscience* doi:10.1101/2020.09.30.321687.
- Oganian Y, Chang EF, 2019. A speech envelope landmark for syllable encoding in human superior temporal gyrus. *Sci. Adv* 14.
- Park H, Ince RAA, Schyns PG, Thut G, Gross J, 2015. Frontal top-down signals increase coupling of auditory low-frequency oscillations to continuous speech in human listeners. *Curr. Biol* 25, 1649–1653. doi:10.1016/j.cub.2015.04.049. [PubMed: 26028433]
- Pefkou M, Arnal LH, Fontolan L, Giraud AL, 2017. θ -Band and β -band neural activity reflects independent syllable tracking and comprehension of time-compressed speech. *J. Neurosci* 37, 7930–7938. doi:10.1523/JNEUROSCI.2882-16.2017. [PubMed: 28729443]
- Pittman-Polletta BR, Wang Y, Stanley DA, Schroeder CE, Whittington MA, Kopell NJ, 2020. Differential contributions of synaptic and intrinsic inhibitory currents to speech segmentation via flexible phase-locking in neural oscillators (preprint). *Neuroscience* doi:10.1101/2020.01.11.902858.
- Ray S, Crone NE, Niebur E, Franaszczuk PJ, Hsiao SS, 2008. Neural correlates of high-gamma oscillations (60–200hz) in macaque local field potentials and their potential implications in electrocorticography. *J. Neurosci. Off. J. Soc. Neurosci* 28, 11526–11536. doi:10.1523/JNEUROSCI.2848-08.2008.
- Ray S, Maunsell JHR, 2011. Different origins of gamma rhythm and high-gamma activity in macaque visual cortex. *PLoS Biol* 9, e1000610. doi:10.1371/journal.pbio.1000610.g008. [PubMed: 21532743]
- Remington ED, Egger SW, Narain D, Wang J, Jazayeri M, 2018. A dynamical systems perspective on flexible motor timing. *Trends Cogn. Sci* 22, 938–952. doi:10.1016/j.tics.2018.07.010. [PubMed: 30266152]
- Russo AA, Bittner SR, Perkins SM, Seely JS, London BM, Lara AH, Miri A, Marshall NJ, Kohn A, Jessell TM, Abbott LF, Cunningham JP, Churchland MM, 2018. Motor cortex embeds muscle-like commands in an untangled population response. *Neuron* 97, 953–966. doi:10.1016/j.neuron.2018.01.004. e8. [PubMed: 29398358]
- Russo AA, Khajeh R, Bittner SR, Perkins SM, Cunningham JP, Abbott LF, Churchland MM, 2020. Neural trajectories in the supplementary motor area and motor cortex exhibit distinct geometries, compatible with different classes of computation. *Neuron* 107, 745–758. doi:10.1016/j.neuron.2020.05.020. e6. [PubMed: 32516573]
- Scheffer-Teixeira R, Belchior H, Leão RN, Ribeiro S, Tort ABL, 2013. On high-frequency field oscillations (>100Hz) and the spectral leakage of spiking activity. *J. Neurosci* 33, 1535–1539. doi:10.1523/JNEUROSCI.4217-12.2013. [PubMed: 23345227]

- Seely JS, Kaufman MT, Ryu SI, Shenoy KV, Cunningham JP, Churchland MM, 2016. Tensor analysis reveals distinct population structure that parallels the different computational roles of areas m1 and V1. *PLOS Comput. Biol* 12, e1005164. doi:10.1371/journal.pcbi.1005164. [PubMed: 27814353]
- Stringer C, Pachitariu M, Steinmetz N, Carandini M, Harris KD, 2019. High-dimensional geometry of population responses in visual cortex. *Nature* 571, 361–365. doi:10.1038/s41586-019-1346-5. [PubMed: 31243367]
- Suzuki M, Larkum ME, 2017. Dendritic calcium spikes are clearly detectable at the cortical surface. *Nat. Commun* 8, 276. doi:10.1038/s41467-017-00282-4. [PubMed: 28819259]
- Theunissen FE, David SV, Singh NC, Hsu A, Vinje WE, Gallant JL, 2001. Estimating spatio-temporal receptive fields of auditory and visual neurons from their responses to natural stimuli. *Netw. Comput. Neural Syst* 12, 289–316. doi:10.1080/net.12.3.289.316.
- Vyas S, Golub MD, Sussillo D, Shenoy KV, 2020. Computation through neural population dynamics. *Annu. Rev. Neurosci* 43, 249–275. doi:10.1146/annurev-neuro-092619-094115. [PubMed: 32640928]
- Wang J, Narain D, Hosseini EA, Jazayeri M, 2018. Flexible timing by temporal scaling of cortical responses. *Nat. Neurosci* 21, 102–110. doi:10.1038/s41593-017-0028-6. [PubMed: 29203897]
- Yi HG, Leonard MK, Chang EF, 2019. The encoding of speech sounds in the superior temporal gyrus. *Neuron* 102, 1096–1110. doi:10.1016/j.neuron.2019.04.023. [PubMed: 31220442]
- Yoon TJ, Cole J, Hasegawa-Johnson M, 2007. August. On the edge: Acoustic cues to layered prosodic domains. *Proceedings of ICPhS* 16, 1264–1267.

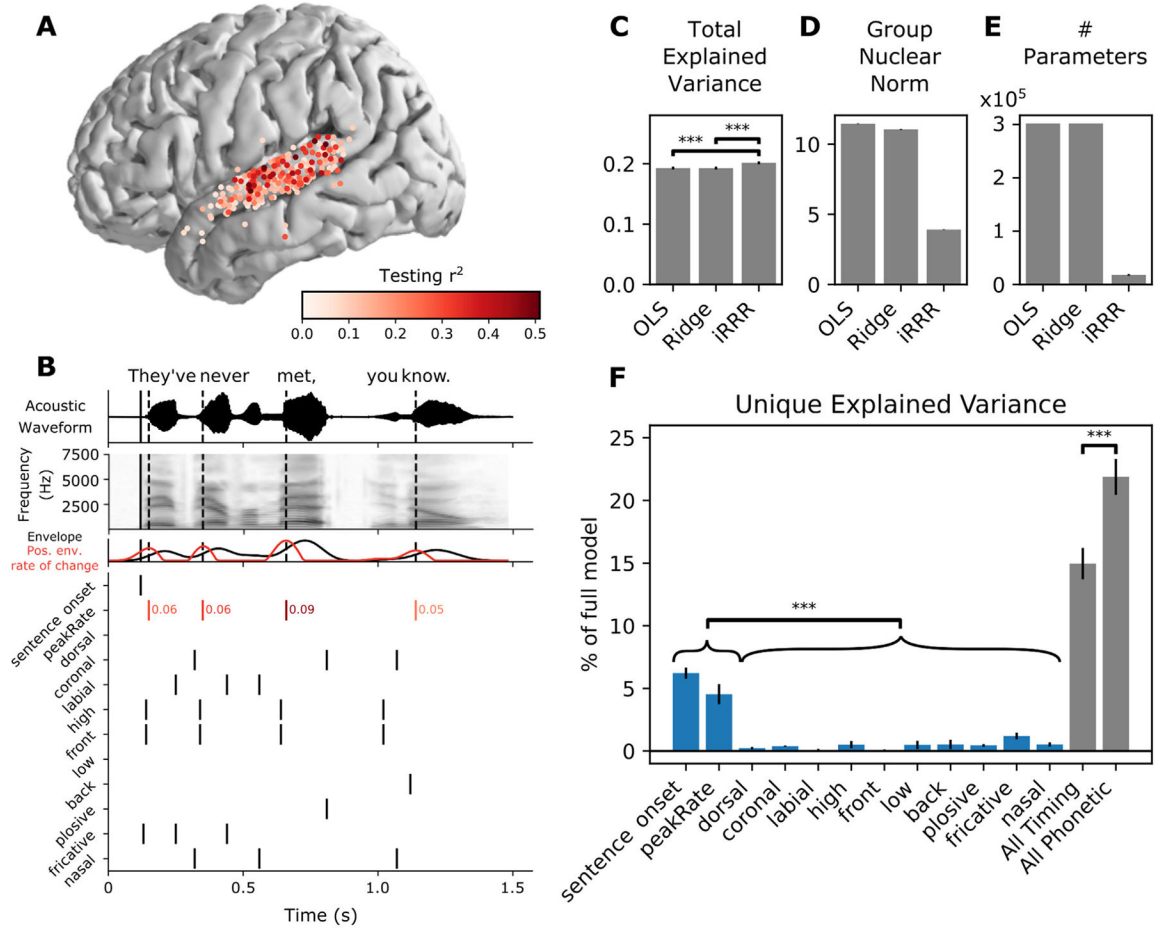
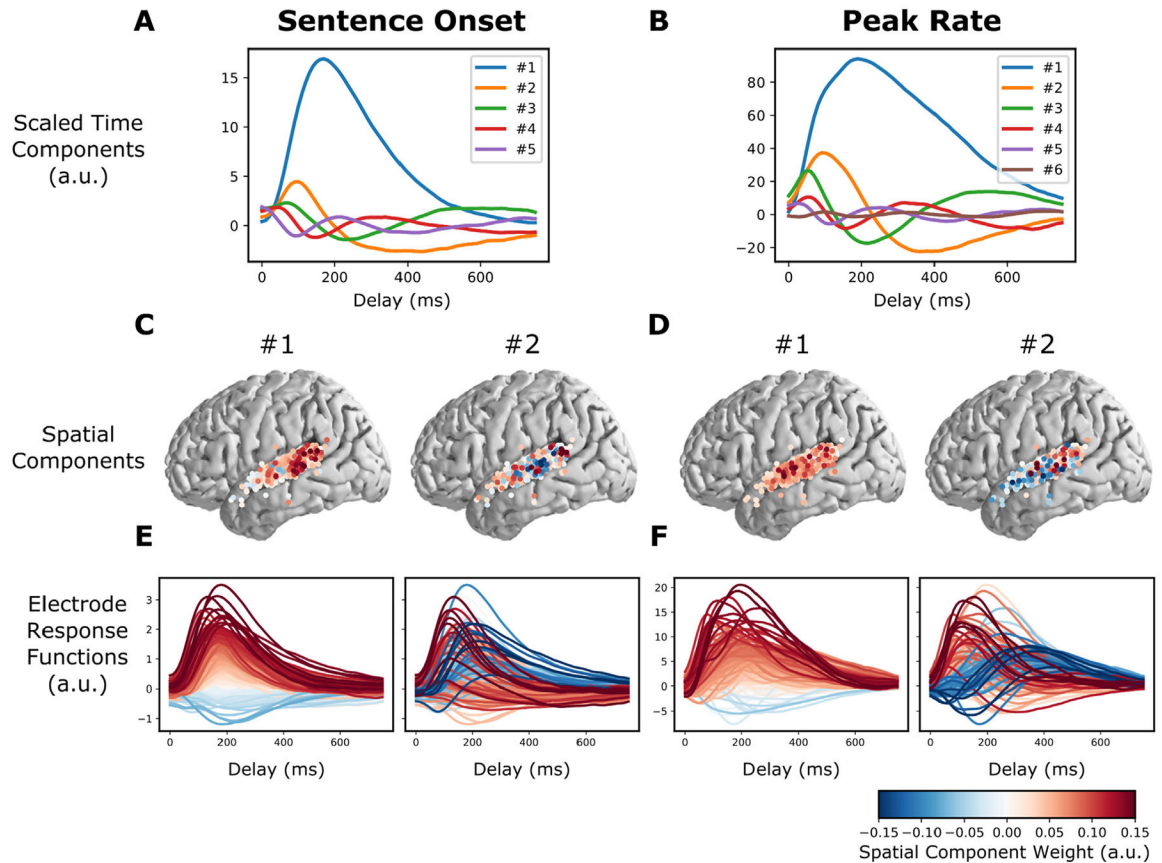


Fig. 1. iRRR outperforms models that treat each electrode individually, and sentence onset and peak rate capture more of the variance than phonetic features. **A:** Electrodes used for model fitting, colored according to the testing r^2 of the linear spectrotemporal (STRF) model (electrodes were selected for subsequent analysis if they were located over STG and if their testing r^2 for the spectrotemporal model was greater than 0.05). **B:** Features used for feature temporal receptive field modeling. Top: the acoustic waveform of an example sentence. The solid vertical line shows the sentence onset event, and the dashed vertical lines show the times of the peak rate events. Second panel: the corresponding mel-band spectrogram. Third panel: the envelope of the acoustic waveform (black) and the positive rate of change of the envelope (red). The peaks in the positive envelope rate of change are the peak rate events. Bottom: the feature time series. White space represents no event (encoded by 0 in the feature matrix), black lines represent event times (encoded by 1), and red lines indicate peak rate event times with their corresponding magnitude indicated to the right. **C, D, E:** Performance of the iRRR model in comparison to ordinary least squares (OLS) and ridge regression (Ridge). 95% confidence intervals were estimated using the standard error of the mean across cross-validation folds (see Section 3.8). Significance was assessed for comparisons using two-sided paired t-tests across cross-validation folds, *** $p < 0.0005$. **C:** Total explained variance, computed as the testing r^2 computed over all speech-responsive

electrodes. D: Group nuclear norm, meaning the penalty term from the iRRR model (see Eq. (11)). E: The effective number of parameters for the fitted models. F: Unique explained variance for each feature (over all speech-responsive electrodes), expressed as a percentage of the variance captured by the full model. Comparing individual features, both timing features have significantly more unique explained variance than all phonetic features, after Bonferroni correction over pairs (left). Also shown is the unique explained variance for the combined timing features (sentence onset and peak rate) and the combined phonetic features (right). When the features are grouped, the phonetic features capture more unique explained variance than the timing features.

**Fig. 2.**

The model fit captures known response differences between pSTG and mSTG. A and B: Time components for the sentence onset and peak rate response matrices, scaled by their singular value (all panels of this figure use the fit from the first cross-validation fold). C: The first two spatial components (across electrodes) for sentence onset. E: The electrode responses to sentence onset events (rows of the sentence onset response matrix), colored by the first (left) or second (right) peak rate spatial component. The first spatial component for sentence onset shows that electrodes with large sentence onset responses (red lines in the left plot of E) tend to be in posterior STG (red circles in the left plot of C). D and F: (like C and E, but for peak rate). The second spatial component divides electrodes into fast and slow peak rate responses (red and blue lines in the right plot of F), which tend to occur over pSTG and mSTG, respectively (red and blue circles in the right plot of D).

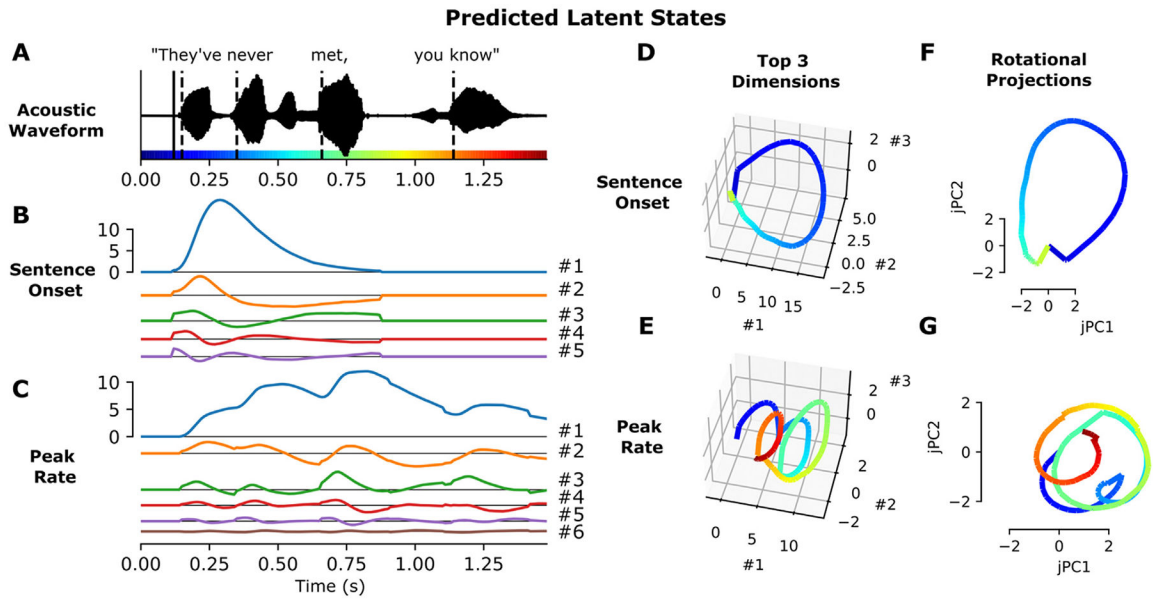


Fig. 3. Feature latent states have rotational dynamics that capture continuous relative timing information. **A:** Acoustic waveform of the stimulus. Solid and dashed vertical lines indicate the timing of the sentence onset and peak rate events, respectively. Colors along the x-axis are used to indicate time in parts **D-G**. **B, C:** Predicted latent states for the sentence onset and peak rate features corresponding to the given stimulus. **D, E:** Top three dimensions of the predicted sentence onset and peak rate latent states (the top three dimensions capture 98.7% and 98.8% of the variance in the sentence onset and peak rate coefficient matrices, respectively). **F, G:** Projection of the predicted sentence onset and peak rate latent states onto the plane of fastest rotation (identified using jPCA). The displayed jPCA projections capture 31.8% and 20.3% of the variance in the sentence onset and peak rate coefficient matrices, respectively. All panels of this figure use the fit from the first cross-validation fold.

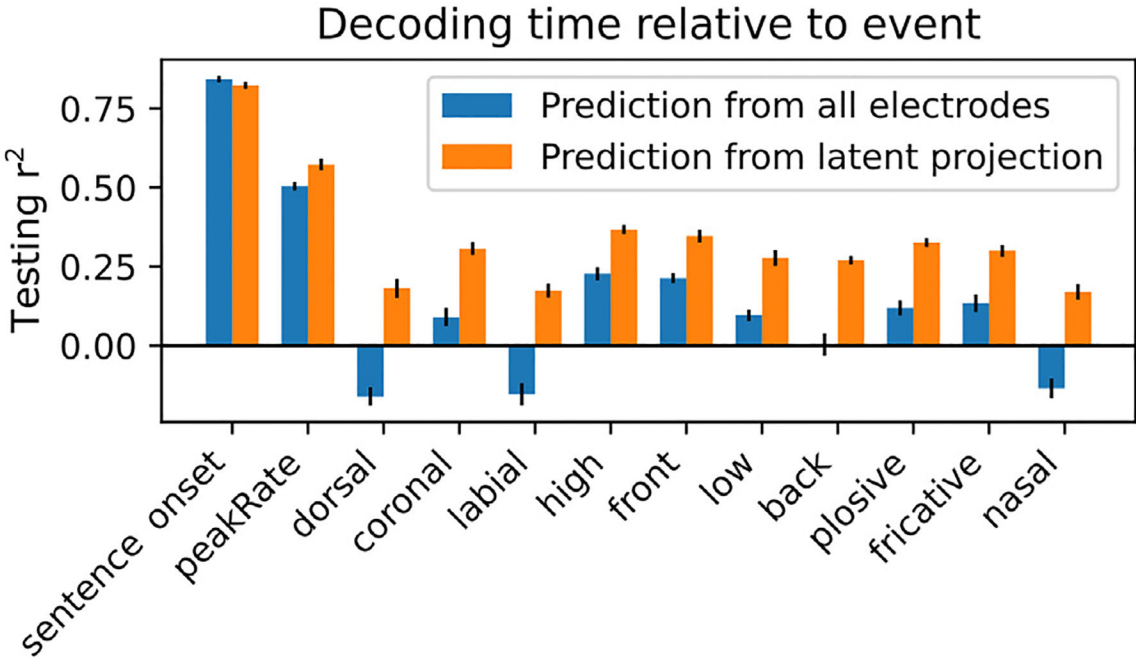


Fig. 4. Latent states from the model can be used to decode time relative to feature events. Performance of a perceptron model trained to decode the time relative to the most recent feature event, for each feature. The models were trained either using the full high-dimensional set of high gamma responses across electrodes (blue bars) or using the projection of those responses onto the subspaces spanned by the feature latent states (orange bars). Performance is quantified using the testing set r^2 .

Author Manuscript

Author Manuscript

Author Manuscript

Author Manuscript

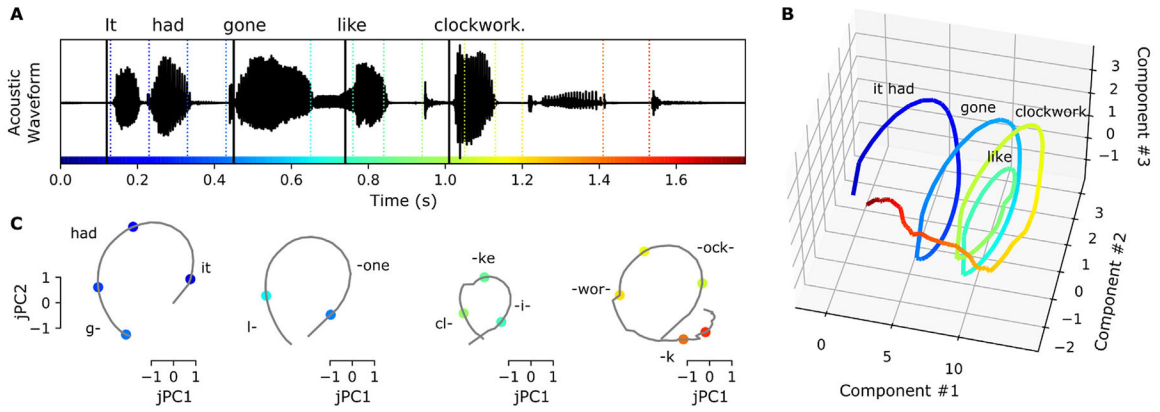


Fig. 5. Peak rate rotational latent states could provide a temporal scaffolding on which individual acoustic features can be organized. A: The acoustic waveform for the stimulus “It had gone like clockwork”. Solid vertical lines indicate the times of peak rate events, and colored dashed vertical lines indicate the times of phonetic feature events. Colors are used to indicate time in all panels. B: The predicted peak rate latent state follows a spiral trajectory in the top 3 dimensions. C: Projected onto the plane of greatest rotation (jPC1 and 2), the predicted peak rate latent state divides the sentence into four intervals, each consisting of a rotation through state space that captures the time since the peak rate event occurred. Downstream processing could combine the relative time information encoded in the peak rate subspace (grey traces) with the feature identities encoded in the feature subspaces (colored points) to compose higher-order representations of words or small groups of words. Text in panels B and C indicates the approximate timing of the words in the stimulus.