RESEARCH ARTICLE

# Spectrally specific temporal analyses of spike-train responses to complex sounds: A unifying framework

**Satyabrata Parida**[1], **Hari Bharadwaj**[1,2], **Michael G. Heinz**[1,2]*

**1** Weldon School of Biomedical Engineering, Purdue University, West Lafayette, Indiana, United States of America, **2** Department of Speech, Language, and Hearing Sciences, Purdue University, West Lafayette, Indiana, United States of America

* mheinz@purdue.edu

## Abstract

Significant scientific and translational questions remain in auditory neuroscience surrounding the neural correlates of perception. Relating perceptual and neural data collected from humans can be useful; however, human-based neural data are typically limited to evoked far-field responses, which lack anatomical and physiological specificity. Laboratory-controlled preclinical animal models offer the advantage of comparing single-unit and evoked responses from the same animals. This ability provides opportunities to develop invaluable insight into proper interpretations of evoked responses, which benefits both basic-science studies of neural mechanisms and translational applications, e.g., diagnostic development. However, these comparisons have been limited by a disconnect between the types of spectrotemporal analyses used with single-unit spike trains and evoked responses, which results because these response types are fundamentally different (point-process versus continuous-valued signals) even though the responses themselves are related. Here, we describe a unifying framework to study temporal coding of complex sounds that allows spike-train and evoked-response data to be analyzed and compared using the same advanced signal-processing techniques. The framework uses a set of peristimulus-time histograms computed from single-unit spike trains in response to polarity-alternating stimuli to allow advanced spectral analyses of both slow (envelope) and rapid (temporal fine structure) response components. Demonstrated benefits include: (1) novel spectrally specific temporal-coding measures that are less confounded by distortions due to hair-cell transduction, synaptic rectification, and neural stochasticity compared to previous metrics, e.g., the correlogram peak-height, (2) spectrally specific analyses of spike-train modulation coding (magnitude and phase), which can be directly compared to modern perceptually based models of speech intelligibility (e.g., that depend on modulation filter banks), and (3) superior spectral resolution in analyzing the neural representation of nonstationary sounds, such as speech and music. This unifying framework significantly expands the potential of preclinical animal models to advance our understanding of the physiological correlates of perceptual deficits in real-world listening following sensorineural hearing loss.

## Author summary

Despite major technological and computational advances, we remain unable to match human auditory perception using machines, or to restore normal-hearing communication for those with sensorineural hearing loss. An overarching reason for these limitations is that the neural correlates of auditory perception, particularly for complex everyday sounds, remain largely unknown. Although neural responses can be measured in humans noninvasively and compared with perception, these evoked responses lack the anatomical and physiological specificity required to reveal underlying neural mechanisms. Single-unit spike-train responses can be measured from preclinical animal models with well-specified pathology; however, the disparate response types (point-process versus continuous-valued signals) have limited application of the same advanced signal-processing analyses to single-unit and evoked responses required for direct comparison. Here, we fill this gap with a unifying framework for analyzing both spike-train and evoked neural responses using advanced spectral analyses of both the slow and rapid response components that are known to be perceptually relevant for speech and music, particularly in challenging listening environments. Numerous benefits of this framework are demonstrated here, which support its potential to advance the translation of spike-train data from animal models to improve clinical diagnostics and technological development for real-world listening.

This is a *PLOS Computational Biology* Methods paper.

## Introduction

Normal-hearing listeners demonstrate excellent acuity while communicating in complex environments. In contrast, hearing-impaired listeners often struggle in noisy situations, even with state-of-the-art intervention strategies (e.g., digital hearing aids). In addition to improving our understanding of the auditory system, the clinical outcomes of these strategies can be improved by studying how the neural representation of complex sounds relates to perception in normal and impaired hearing. Numerous electrophysiological studies have explored the neural representation of perceptually relevant sounds in humans using evoked far-field recordings, such as frequency following responses (FFRs) and electroencephalograms [1–3]. Note that we use *electrophysiology* and *neurophysiology* to refer to evoked far-field responses and single-unit responses, respectively (see S1 Table for glossary). While these evoked responses are attractive because of their clinical viability, they lack anatomical and physiological specificity. Moreover, the underlying sensorineural hearing loss pathophysiology is typically uncertain in humans. In contrast, laboratory-controlled animal models of various pathologies can provide specific neural correlates of perceptual deficits that humans experience, and thus hold great scientific and translational (e.g., pharmacological) potential. In order to synergize the benefits of both these approaches to advance basic-science and translational applications to real-world listening, two major limitations need to be addressed.

First, there exists a significant gap in relating spike-train data recorded invasively from animals and evoked noninvasive far-field recordings feasible in humans (and animals) because

the two signals are fundamentally different in form (i.e., binary-valued point-process data versus continuous-valued signals). While the continuous nature of the evoked-response amplitude allows for any of the advanced signal-processing techniques developed for continuous-valued signals to be applied (e.g., multitaper approaches to robust spectral estimation [4]), spike-train analyses have been much more limited (e.g., in their application to real-world signals, as reviewed in S1 Text). This is a critical gap because most perceptual deficits and machine-hearing limits occur for speech in noise rather than for speech in quiet [5, 6]. For example, classic neurophysiological studies have quantified the temporal coding of stationary and periodic stimuli using metrics such as vector strength (VS [7–9], also see S1 Appendix), whereas more recent correlogram analyses have provided temporal-coding metrics for non-periodic stimuli, such as noise [10, 11]. However, as reviewed in S1 Text, these metrics can be influenced by distortions from nonlinear cochlear processes [12, 13], and often ignore response phase information that is likely to be perceptually relevant for simple tasks [14] as well as for speech intelligibility [15, 16].

A second important gap exists because current spectrotemporal tools to evaluate temporal coding in the auditory system are largely directed at processing of stationary signals by linear and time-invariant systems. However, the auditory system exhibits an array of nonlinear (e.g., two-tone suppression, compressive gain, and rectification) and time-varying (e.g., adaptation and efferent feedback) mechanisms [17, 18]. These mechanisms interact with nonstationary stimulus features (e.g., frequency transitions and time-varying intensity fluctuations, Fig 1A and 1B) to shape the neural coding and perception of these signals [19–21]. In fact, the response of an auditory-nerve (AN) fiber to even a simple stationary tone shows nonstationary features, such as a sharp onset and adaptation (Fig 1C), illustrating the need for nonstationary analyses of temporal coding. However, the extensive single-unit speech coding studies using classic spike-train metrics have typically been limited to synthesized and stationary speech tokens, which has deferred the study of the rich kinematics present in natural speech [12, 22, 23]. Some windowing-based approaches have been used to study time-varying stimuli and responses [24, 25], but the approaches used have imposed a limit on the temporal and spectral resolution with which dynamics of the auditory system can be studied.

The present study focuses on developing spectrotemporal tools to characterize the neural representation of kinematics naturally present in real-world signals, speech in particular, that are appropriate for the nonlinear and time-varying auditory system. We describe a unifying framework to study temporal coding in the auditory system, which allows direct comparison of single-unit spike-train responses with evoked far-field recordings. In particular, we demonstrate the unifying merit of using alternating-polarity peristimulus time histograms (*apPSTHs*, Table 1), a collection of PSTHs obtained from responses to both positive and negative polarities of the stimulus. By using both polarities, neural coding of natural sounds can be studied using the common temporal dichotomy between the slowly varying envelope (ENV) and rapidly varying temporal fine structure (TFS) (Fig 1E and 1F), which has been especially relevant for speech-perception studies [26, 27]. We derive explicit relations between *apPSTHs* and existing metrics for quantifying temporal coding in auditory neurophysiology (reviewed in S1 Text), namely VS and correlograms, to show that no information is lost by using *apPSTHs*. In fact, the use of *apPSTHs* is computationally more efficient, provides more precise spectral estimators, and opens up new avenues for perceptually relevant analyses that are otherwise not possible. Next, an *apPSTH*-based ENV/TFS taxonomy is presented, including existing and new metrics. This taxonomy allows for spectrally specific analyses that avoid distortions due to inner-hair-cell transduction and synaptic rectification processes, resulting in more accurate characterizations of temporal coding than with previous metrics. Finally, these methods are extended in novel ways to include the study of nonstationary signals at superior
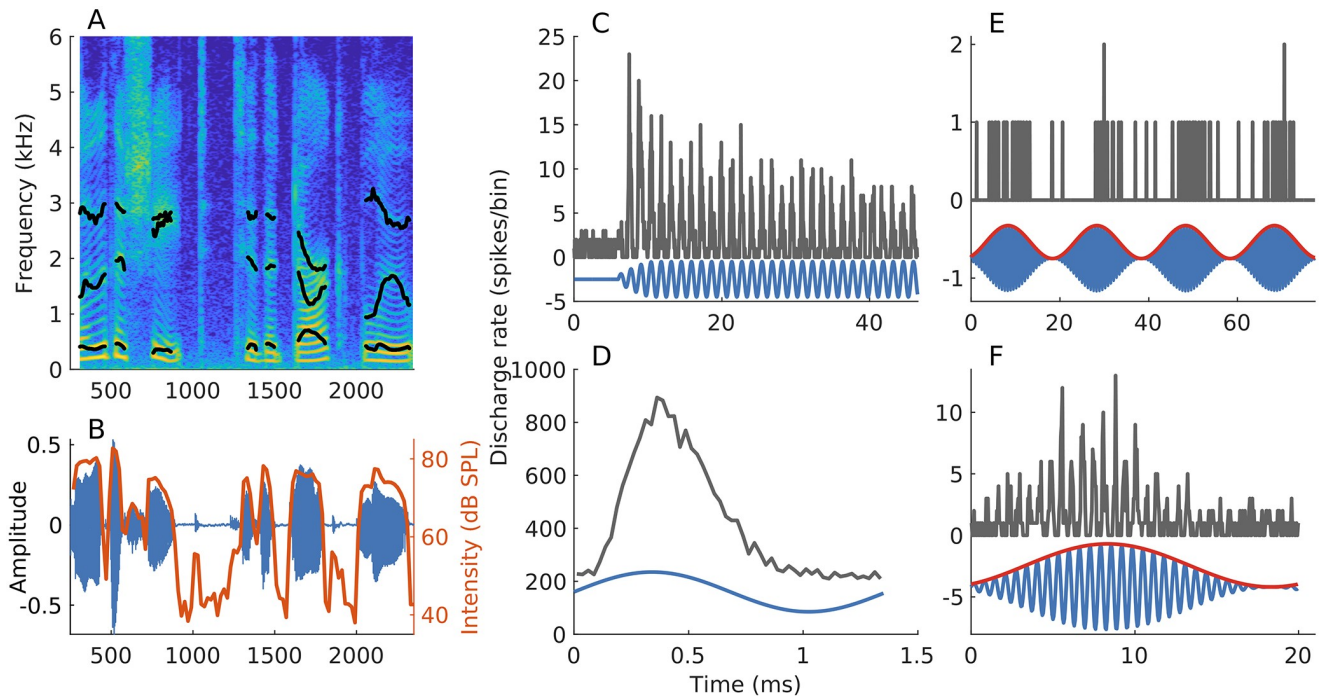
**Fig 1. Neural responses of AN fibers are invariably nonstationary, even when the stimulus is not.** (A, B) Spectrogram and waveform of a speech segment ($s_4$ described in *Materials and Methods*). Formant trajectories (black lines in panel A) and short-term intensity (red line in panel B, computed over 20-ms windows with 80% overlap) vary with time, highlighting two nonstationary aspects of speech stimuli. (C) PSTH constructed using spike trains in response to a tone at the AN-fiber's characteristic frequency (CF, most-sensitive frequency; fiber had CF = 730 Hz and was high spontaneous rate or SR [28]). Tone intensity = 40 dB SPL. Even though the stimulus is stationary, the response is nonstationary (i.e., sharp onset followed by adaptation). (D) Period histogram, constructed from the data used in C, demonstrates the phase-locking ability of neurons to individual stimulus cycles. (E) PSTH constructed using spike trains in response to a sinusoidally amplitude-modulated (SAM) CF-tone (50-Hz modulation frequency, 0-dB modulation depth, 35 dB SPL) from an AN fiber (CF = 1.4 kHz, medium SR). (F) Period histogram (for one modulation period) constructed from the data used in E. The response to the SAM tone follows both the modulator (envelope, red, panels E and F) as well as the carrier (temporal fine structure), the rapid fluctuations in the signal (blue, panel F). Bin width = 0.5 ms for histograms in C-F. Number of stimulus repetitions for C and E were 300 and 16, respectively.

https://doi.org/10.1371/journal.pcbi.1008155.g001

**Table 1. *apPSTH*-taxonomy for ENV & TFS.**

| PSTH name | Notation: (time, frequency) | Definition | ENV and/or TFS representation | Rectifier distortion | Comments |
|---|---|---|---|---|---|
| Positive | $p(t)$, $P(f)$ | Positive polarity | TFS & ENV | Large | |
| Negative | $n(t)$, $N(f)$ | Negative polarity | TFS & ENV | Large | |
| Difference | $d(t)$, $D(f)$ | $\dfrac{p(t) - n(t)}{2}$ | TFS & ENV | Small | Includes both the carrier and sideband components (thus not a clean representation of TFS) |
| Sum | $s(t)$, $S(f)$ | $\dfrac{p(t) + n(t)}{2}$ | ENV | Large | Consistent representation of spectrally specific modulation strength but confounded by rectifier distortion at $2 \times CF$ |
| Analytic | $a(t)$, $A(f)$ | $d(t) + j\mathcal{H}\{d(t)\}$ | TFS & ENV | Small | $\mathcal{H}\{\cdot\}$ is the Hilbert transform operator |
| Hilbert envelope | $e(t)$, $E(f)$ | $\lvert a(t) \rvert / \sqrt{2}$ | ENV | Small | Polarity-sensitive ENV (subject to TFS phase locking) |
| Hilbert phase | $\phi(t)$, $\Phi(f)$ | $\sqrt{2} \times rms[d(t)] \times cos[\angle a(t)]$ | TFS | Small | Carrier TFS (subject to TFS phase locking) |

We define *apPSTHs* as the collection of PSTHs derived using both polarities of the stimulus. The pair of PSTHs, $p(t)$ and $n(t)$, is a sufficient statistic for *apPSTHs* since all other PSTHs in the group can be derived from the two. Alternatively, the pair, $d(t)$ and $s(t)$, is also a sufficient statistic for *apPSTHs*. Each PSTH (e.g., the positive polarity PSTH) can be expressed in the time domain [$p(t)$] or in the frequency domain [$P(f)$]. A graphical illustration for these *apPSTHs* is in S1 Fig.

https://doi.org/10.1371/journal.pcbi.1008155.t001

spectrotemporal resolution compared to conventional windowing-based approaches, like the spectrogram or wavelet analysis.

## Results

### A unified framework for quantifying temporal coding based on alternating-polarity PSTHs (*apPSTHs*)

In this section, we first show that *apPSTHs* can be used to unify classic metrics, e.g., *VS* and correlograms (reviewed in S1 Text), in a computationally efficient manner. Then, we show that *apPSTHs* offer more precise spectral estimates compared to correlograms and allow for perceptually relevant analyses that are not possible with classic metrics.

**apPSTHs permit computationally efficient temporal analyses.** Let us denote the PSTHs in response to the positive and negative polarities of a stimulus as $p(t)$ and $n(t)$, respectively. Then, the *sum PSTH*, $s(t)$, which represents the polarity-tolerant component in the response, is estimated as

$$s(t) = \frac{p(t) + n(t)}{2}. \tag{1}$$

The *difference PSTH*, $d(t)$, which represents the polarity-sensitive component in the response, is estimated as

$$d(t) = \frac{p(t) - n(t)}{2}. \tag{2}$$

The difference PSTH has been previously described as the compound PSTH [29]. Here we use the terms *sum* and *difference* for $s(t)$ and $d(t)$, respectively, for clarity. Compared to the spectra of the single-polarity PSTHs [i.e., of $p(t)$ or $n(t)$], the spectrum of the difference PSTH, $D(f)$, is substantially less confounded by rectifier-distortion artifacts ([23], also see S1 Fig panels B and D). This improvement occurs because even-order distortions, which strongly contribute to these artifacts, are effectively canceled out by subtracting PSTHs for opposite polarities. A second way spectral peaks absent in the stimulus can arise in the $p(t)$-spectrum is because of propagating combination tones of cochlear origin (e.g., distortion products) [30]. Unlike rectifier distortion, which is an artifact of analysis, combination tones are present in the cochlea and can affect perception. As the phase of these combination tones depends on stimulus polarity [30], these perceptually relevant combination tones are captured in the difference PSTH. These distinct sources are discussed in more detail by Young and Sachs with respect to analyses of stationary synthesized-vowel responses from AN fibers [12].

The Fourier magnitude spectrum of the difference PSTH has been referred to as the synchronized rate. We show that the synchronized rate relates to *VS* by

$$VS(f) = \frac{|D(f)|}{N}, \tag{3}$$

where $f$ is frequency in Hz, and $N$ is the total number of spikes (S2 Appendix).

In addition, we demonstrate that the autocorrelogram and the shuffled autocorrelation (SAC) function of the PSTH are related (S3 Appendix), which leads to important computational efficiencies. In particular the SAC for a set of M spike trains $X = \{\underline{x_1}, \underline{x_2}, ..., \underline{x_M}\}$ can be estimated as

$$SAC(X) = \mathcal{R}_{\mathcal{X}}(PSTH_X) - \sum_{i=1}^{M} \mathcal{R}_{\mathcal{X}}(\underline{x_i}), \tag{4}$$

where $\mathcal{R}_{\mathcal{X}}$ is the autocorrelation operator, and $PSTH_X$ is the PSTH constructed using $X$. Similarly, the SCC for two sets of spike trains $X = \{\underline{x_1}, \underline{x_2}, ..., \underline{x_L}\}$ and $Y = \{\underline{y_1}, \underline{y_2}, ..., \underline{y_M}\}$ can be estimated as

$$SCC(X, Y) = \mathcal{R}_{\mathcal{X}\mathcal{Y}}(PSTH_X, PSTH_Y), \qquad (5)$$

where $PSTH_X$ and $PSTH_Y$ are PSTHs constructed using $X$ and $Y$, respectively, and $\mathcal{R}_{\mathcal{X}\mathcal{Y}}$ is the cross-correlation operator. Since SACs and SCCs can be computed using *apPSTHs*, it follows that *sumcor* and *difcor* can also be computed using *apPSTHs* (S4 Appendix). As *apPSTHs* can be used to compute correlograms, *apPSTHs* offer the same degree of smoothing as correlograms.

Importantly, the use of *apPSTHs* to compute correlograms is computationally more efficient compared to the existing correlogram-estimation method, i.e., by tallying all interspike intervals. For a fixed stimulus duration and PSTH resolution, estimating the autocorrelation function of the PSTH requires constant time complexity [$\mathcal{O}(1)$]. Thus, for $N$ spikes, the SAC and SCC can be computed with $\mathcal{O}(N)$ complexity that is needed for constructing the PSTH using Eqs 4 and 5. This is substantially better than the $\mathcal{O}(N^2)$ complexity needed to compute the correlograms by tallying shuffled all-order interspike intervals. For example, consider a spike-train dataset that consists of 50 repetitions of a stimulus with 100 spikes per repetition. To compute the SAC using (all-order) ISIs, each spike time (5000 unique spikes) has to be compared with spike times from all other repetitions (4900 spike times). This tallying method requires $24.5 \times 10^6$ (i.e., $5000 \times 4900$) operations to compute the SAC, where one operation consists of comparing two spike times and incrementing the corresponding SAC-bin by 1. In contrast, only 5000 operations are needed to construct the PSTH for 5000 ($50 \times 100$) total spikes. The PSTH can then be used to estimate the SAC with constant time complexity. In addition to their computational efficiency, *apPSTHs* offer additional benefits for relating single-unit responses to far-field responses, for spectral estimation, and for speech-intelligibility modeling, as discussed below.

**apPSTHs unify single-unit and far-field analyses.** The PSTH is particularly attractive because the PSTH from single neurons or a population of neurons, by virtue of being a continuous signal, can be directly compared to evoked potentials in response to the same stimulus (e.g., Fig 2). In this example, the speech sentence $s_3$ was used to record the frequency following response (FFR) from one animal. The same stimulus was also used to record spike trains from AN fibers (N = 246) from 13 animals. The mean $d(t)$ and mean $s(t)$ were computed by pooling PSTHs across all neurons. The difference and sum FFRs were estimated by subtracting and averaging FFRs to alternating polarities, respectively. This approach of estimating polarity-tolerant and polarity-sensitive FFR components is well established [31–33]. Qualitatively, the periodicity information in the mean $d(t)$ and the difference FFR were similar (Fig 2A); this is expected because the difference FFR receives significant contributions from the auditory nerve [34]. To compare the spectra for the two responses, a 100-ms segment was considered. The first formant ($F_1$) and the first few harmonics of the fundamental frequency ($F_0$) were well captured in both spectra. $F_2$ was also well captured in the difference FFR, and to a lesser extent, in the mean $d(t)$.

The mean $s(t)$ and the sum FFR also show comparable temporal features in these nonstationary responses (Fig 2C). For example, both responses show sharp onsets for plosive and fricative consonants. The segment considered in Fig 2B was used to compare the spectra for the two sum responses. Both spectra show similar spectral peaks near the first two harmonics of $F_0$ (Fig 2D), which indicates that pitch-related periodicity is well captured in both the sum FFR and mean $s(t)$. However, there are some discrepancies between the relative heights of the
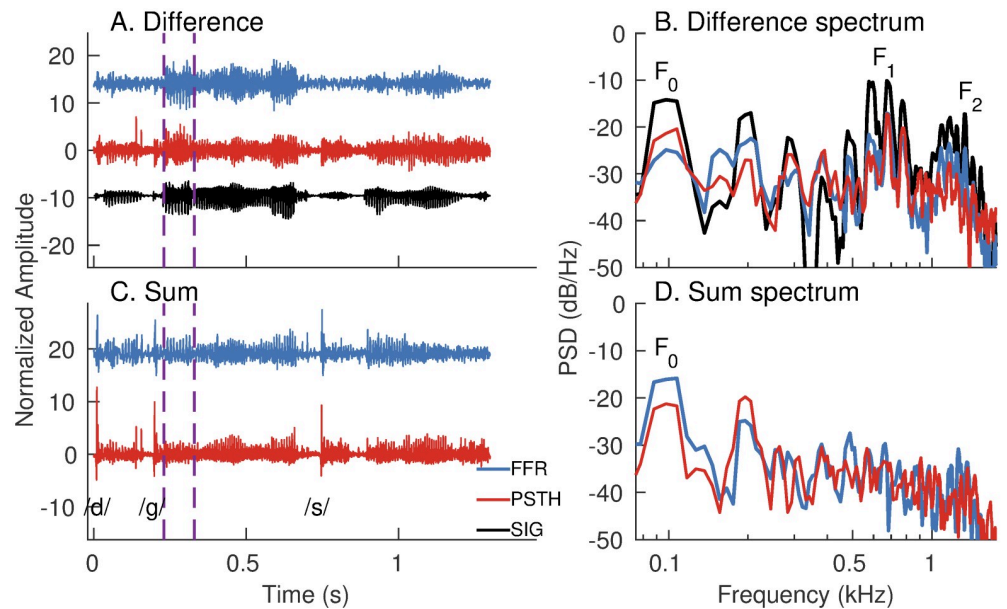
**Fig 2. *apPSTHs* can be directly compared to evoked potentials in response to the same stimulus.** (A) Time-domain waveforms for the difference FFR (blue) and mean difference PSTH [$d(t)$, red] in response to a Danish speech stimulus, $s_3$ (black). Mean $d(t)$ was computed by taking the grand average of $d(t)s$ from 246 AN fibers from 13 animals (CFs: 0.2 to 11 kHz). The difference FFR was estimated by subtracting FFRs to alternating stimulus polarities. (B) Spectra for the signals in A for a 100-ms segment (purple dashed lines in A). (C) Time-domain waveforms for the sum FFR (blue) and mean sum PSTH [$s(t)$, red] for the same stimulus. Both responses show sharp onsets for plosive (/d/ and /g/) and fricative (/s/) consonants. (D) Spectra for the responses in C for the same segment considered in B. The mean $s(t)$ was estimated as the grand average of $s(t)s$ from 246 neurons. Sum FFR was estimated by halving the sum of the FFRs to both polarities. Stimulus intensity = 65 dB SPL.

first two $F_0$-harmonics. These could arise because the average FFR primarily reflects activity of high-frequency neurons from rostral generators (e.g., the inferior colliculus) [34], which show stronger polarity-tolerant responses compared to the auditory nerve [35]. In contrast, the mean $s(t)$ is based on responses of AN fibers, which show strong polarity-sensitive responses to $F_0$ due to tuning-curve tail responses at high sound levels like that used here. These tail responses contribute to power at $2F_0$ as rectifier distortion. Other potential sources that can contribute to any far-field evoked response include receptor potentials (e.g., cochlear microphonic, or the CM) and electrical interference. However, CM is substantially reduced in summed mass responses (since odd harmonics cancel), although CM may not be completely removed because even harmonics remain [36, 37]. In fact, destructive interference between the auditory neurophonic and the CM has been seen previously in mass responses [37], which could reduce the $2F_0$ component in the sum-FFR spectrum but not in the *apPSTH* for which the CM is not present (Fig 2D). Electrical interference had insignificant effect on these FFR data (Fig 2 in [38]). In general, these sources can substantially contribute to evoked responses, such as the compound action potential, and thus should be considered when comparing these evoked responses with invasive spike-train data [37]. In this regard, using the *apPSTH*-based framework to analyze invasive spike-train recordings allows direct comparison of invasive single-unit data with noninvasive continuous-valued evoked potentials and evaluation of the neural origins of evoked responses.

**Variance of *apPSTH*-based spectral estimates can be reduced relative to correlogram-based spectral estimates.** Temporal information in a signal can be studied not only in the

time domain (e.g., using correlograms) but also in the frequency domain (e.g., using the power spectral density, PSD). The frequency-domain representation often provides a compact alternative compared to the time-domain counterpart. In the framework of spectral estimation, the source ("true") spectrum, which is unknown, is regarded as a parameter of a random process that is to be estimated from the available data (i.e., from examples of the random process). Spectral estimation is complicated by two factors: (1) finite response length, and (2) stochasticity of the system. The former introduces bias to the estimate, i.e., the PSD at a given frequency can differ from the true value. This bias reflects the leakage due to power at nearby (narrowband bias) and far-away (broadband bias) frequencies (due to the inherent temporal windowing from the finite-duration response). Stochasticity of the system adds randomness to the sampled data, which creates variance in the estimate. Desirable properties of PSD estimators are minimized bias and variance. Bias can be reduced by multiplying the data (prior to spectral estimation) with a taper that has a strong energy concentration near 0 Hz. Variance can be reduced by using a greater number of tapers to estimate multiple (independent) PSD estimates, which can be averaged to compute the final estimate. The multitaper approach optimally reduces the bias and variance of the PSD estimate [4, 39]. In this approach, for a given data length, a frequency resolution is chosen, based on which a set of orthogonal tapers are computed. These tapers include both even and odd tapers, which can be used to obtain the independent PSD estimates to be averaged. In contrast, for the same frequency resolution, only even tapers can be used with correlograms as they are even sequences [40, 41]. Therefore, variance in the PSD estimate can be reduced by a factor of up to 2 by using *apPSTHs* instead of correlograms.

For example, the benefit (in terms of spectral-estimation variance) of using the multitaper spectrum of $d(t)$, as opposed to the common approach of estimating the discrete Fourier transform (DFT) of the *difcor*, can be quantified by comparing the two spectra at a single frequency (Fig 3). Here, a 100-ms segment of the $s_3$ speech stimulus was used as the analysis window. The segment had an $F_0$ of 98 Hz and $F_1$ of 630 Hz (Fig 3A). Fig 3B shows example spectra estimated using spike trains recorded from a low-frequency AN fiber [CF = 900 Hz, SR = 81 spikes/s]. The multitaper spectrum was estimated using the MATLAB function *pmtm* (two tapers corresponding to a time-bandwidth product of 3, adaptive weights [4]). To compare variances in the two estimated spectra, fractional power at the 6th harmonic was considered, as this harmonic was closest to $F_1$. This analysis was restricted to neurons (N = 10) for which data was available for at least 75 repetitions per polarity and that had a CF between 0.3 and 2 kHz. For each neuron, 25 spike trains per polarity were chosen randomly 12 times to estimate fractional power at the 6th harmonic. The same set of spike trains were used to estimate distributions for both the *difcor*-spectrum and $D(f)$. The ratio of *difcor*-based fractional power variance to the *apPSTH*-based fractional power variance at $6F_0$ was >1 for all 10 neurons considered (Fig 3D), demonstrating the benefit of being able to compute a multitaper spectrum from $d(t)$ compared to the *difcor*-spectrum in reducing variance. Overall, these results indicate that less data are required to achieve the same level of precision in a spectral metric based on the multitaper spectrum of an *apPSTH* compared to the same metric derived from the DFT of the correlogram.

**Benefits of *apPSTHs* for speech-intelligibility modeling.** Speech-intelligibility (SI) models aim to predict the effects of acoustic manipulations of speech on perception. Thus, SI models allow for quantitative evaluation of the perceptually relevant features in speech. More importantly, SI models can guide the development of optimal hearing-aid strategies for hearing-impaired listeners. However, state-of-the-art SI models are largely based on the acoustic signal, where there is no physiological basis to capture the various effects of sensorineural hearing loss (SNHL) [16, 42–45]. In contrast, neurophysiological SI models (i.e., SI models based
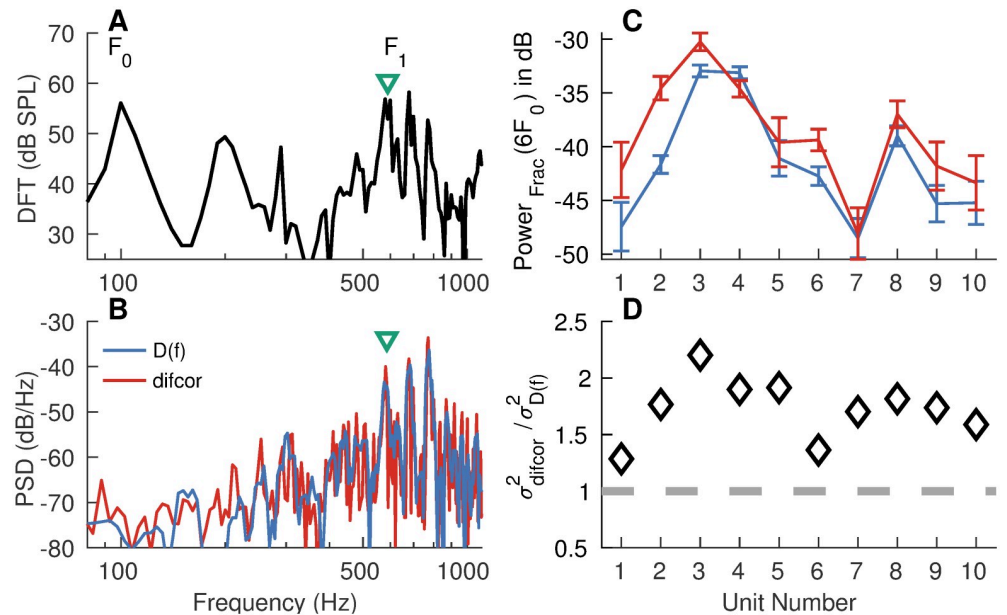
**Fig 3. Lower spectral-estimation variance can be achieved using *apPSTHs* (with multiple tapers) compared with *difcor* correlograms.** (A) Spectrum for the 100-ms segment in the speech sentence *s3* ($F_0 \sim 98$ Hz, $F_1 \sim 630$ Hz) used for analysis. (B) Example spectra for an AN fiber (CF = 900 Hz, high SR) with spikes from 25 randomly chosen repetitions per polarity. The first two discrete-prolate spheroidal sequences were used as tapers corresponding to a time-bandwidth product of 3 to estimate $D(f)$, the spectrum of $d(t)$. No taper (i.e., rectangular window) was used to estimate the *difcor* spectrum. The AN fiber responded to the 6th, 7th and 8th harmonic of the fundamental frequency. (C) Error-bar plots for fractional power (*Power_{Frac}*) at the frequency (green triangle) closest to the 6th harmonic. Error bars were computed for 12 randomly and independently drawn sets of 25 repetitions per polarity. The same spikes were used to compute the spectra for $d(t)$ (blue) and *difcor* (red). (D) Diamonds denote the ratio of variances for the *difcor*-based estimate to the $d(t)$-based estimate. This ratio was greater than 1 (i.e., above the dashed gray line) for all units considered, which demonstrates that the variance for the multitaper-$d(t)$ spectrum was lower than the *difcor*-spectrum variance. AN fibers with CFs between 0.3 and 2 kHz and with at least 75 repetitions per polarity of the stimulus were considered. Bin width = 0.1 ms for PSTHs. Sampling frequency = 10 kHz for FFRs. Stimulus intensity = 65 dB SPL.

on neural data) are particularly important in this regard since spike-train data from preclinical animal models of various forms of SNHL provide a direct way to evaluate the effects of SNHL on speech-intelligibility modeling outcomes [46, 47].

A major advantage of PSTH-based approaches over correlogram-based approaches is that they can be used to extend a wider variety of acoustic SI models to include neurophysiological data. In particular, correlograms can be used to extend power-spectrum-based SI models [42–45, 48] but not for the more recent SI models that require phase information of the response [16, 49]. For example, the speech envelope-power-spectrum model (sEPSM) has been evaluated using simulated spike trains since sEPSM only requires power in the response envelope, which can be estimated from the *sumcor* spectrum [47]. However, *sumcor* cannot be used to evaluate envelope-phase-based SI models since it discards phase information. Studies have shown that the response phase can be important for speech intelligibility [15, 50]. In contrast to the *sumcor*, the time-varying PSTH contains both phase and magnitude information, and thus, can be used to evaluate both power-spectrum- and phase-spectrum-based SI models. For example, because the PSTH $p(t)$ [or $n(t)$] is already rectified, it can be filtered through a modulation filter bank to estimate "internal representations" in the modulation domain (Fig 4). These spike-train-derived "internal representations" are analogous to those used in phase-
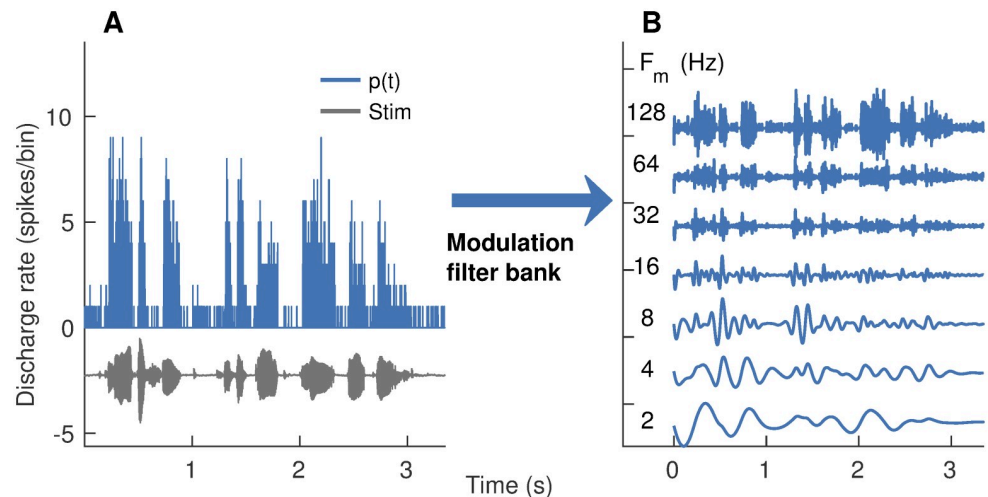
**Fig 4. Modulation-domain internal representations for speech coding can be obtained from PSTH-based envelopes.** PSTH response [$p(t)$] from one AN fiber (CF = 290 Hz, SR = 12 spikes/s) is shown. (A) Time-domain waveforms for the stimulus (gray) and $p(t)$ (blue). (B) Output of a modulation filter bank after the processing of $p(t)$. Modulation filters were zero-phase, fourth-order, and octave-wide IIR filters. Center frequencies ($F_m$) for these filters ranged from 2 to 128 Hz (octave spacing), similar to those used in recent psychophysically based SI models (e.g., [16]). PSTH bin width = 0.5 ms. 15 stimulus repetitions. Stimulus intensity = 60 dB SPL.

spectrum-based SI models [16, 49] and can be further processed by existing SI back-ends to estimate SI values. This example demonstrates a proof of concept of using spike-train data to evaluate a spectrally specific envelope-based SI model using *apPSTHs*. In general, SI models that include a peripheral or modulation filter bank representation, which is the case for most successful SI models (e.g., the speech transmission index [51], the spectrotemporal modulation index [52], speech envelope power spectrum models [48, 53]), can be evaluated using spike-train data recorded from peripheral (e.g., auditory-nerve fibers) or central (e.g., inferior colliculus) neurons, respectively, using *apPSTHs*. Therefore, these analyses allow for the evaluation of a wider variety of acoustic-based SI models in the neural domain (magnitude and phase), where translationally relevant data can be obtained from preclinical animal models of various forms of SNHL.

## Quantifying ENV and TFS using *apPSTHs* for stationary signals

In this subsection, we first describe existing and novel ENV and TFS components that can be derived from *apPSTHs*. Next, we compare relative merits of the novel components over existing ENV and TFS components using simulated data. Finally, we apply *apPSTHs* to analyze spike-train data recorded to speech and speech-like stimuli.

   **Several ENV and TFS components can be derived from *apPSTHs* with spectral specificity.**   The neural response envelope can be obtained from *apPSTHs* in two orthogonal ways: (1) the low-frequency signal, $s(t)$, and (2) the Hilbert envelope of the high-frequency carrier-related energy in $d(t)$. $s(t)$ is thought to represent the polarity-tolerant response component, which has been defined as the envelope response [10, 35]. For a stimulus with harmonic spectrum, $s(t)$ captures the envelope related to the beating between harmonics. In addition, onset and offset responses (e.g., in response to high-frequency fricatives, Fig 2C) are also well captured in $s(t)$. Although *sumcor* and $s(t)$ are related, dynamic features like onset and offset responses are captured in $s(t)$, but not in the *sumcor* since the *sumcor* discards phase

information by essentially averaging ENV coding across the whole stimulus duration. The use of the sum envelope is popular in far-field responses [31–33] but not directly in auditory neurophysiology studies. A major disadvantage of $s(t)$ is that it is affected by rectifier distortions if a neuron phase locks to low-frequency energy in the stimulus (e.g., Fig 5A; discussed further below).
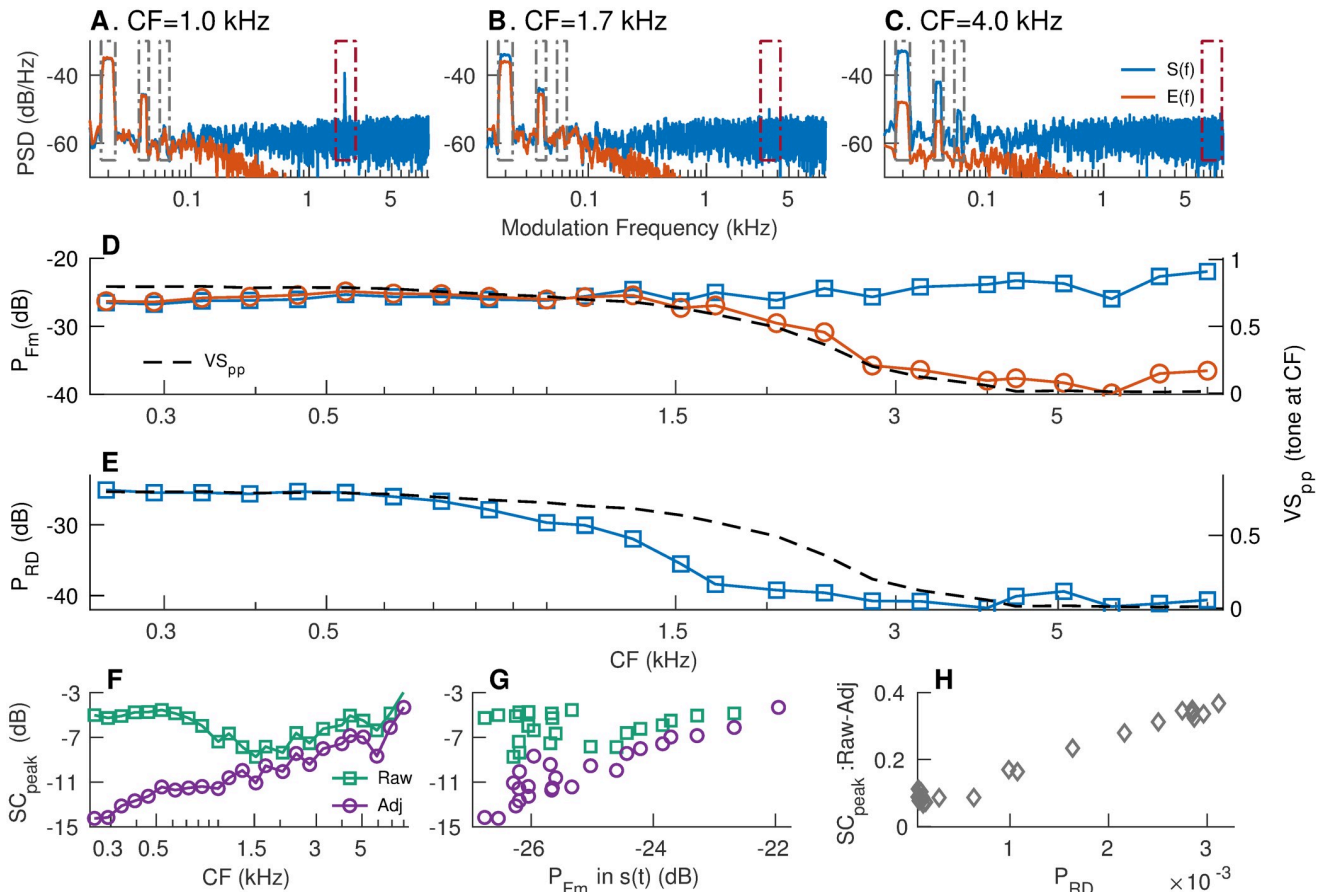


**Fig 5. Envelope-coding metrics should be spectrally specific to avoid artifacts due to rectifier distortion and neural stochasticity.** Simulated responses for 24 AN fibers (log-spaced between 250 Hz and 8 kHz) were obtained using a computational model (parameters listed in S2 Table) using SAM tones at CF (modulation frequency, $F_m$ = 20 Hz; 0-dB (100%) modulation depth) as stimuli. Stimulus intensity ∼ 65 dB SPL. $S(f)$ (blue) and $E(f)$ (red) for three example model fibers with CFs = 1.0, 1.7, and 4 kHz (panels A-C) illustrate the relative merits of $s(t)$ and $e(t)$, and the potential for rectifier distortion to corrupt envelope coding metrics. $d(t)$ was band-limited to a 200-Hz band near $F_c$ for each fiber prior to estimating $e(t)$ from the Hilbert transform of $d(t)$. (A) For the 1-kHz fiber, $S(f)$ and $E(f)$ are nearly identical in the $F_m$ band. $S(f)$ is substantially affected by rectifier distortion at $2 \times CF$, which can be ignored using spectrally specific analyses. (B) The two envelope spectra are largely similar near the $F_m$ bands since phase-locking near the carrier (1.7 kHz) is still strong (panel D). Rectifier distortion in $S(f)$ is greatly reduced since phase-locking at twice the carrier frequency (3.4 kHz) is weak. (C) $F_m$-related power in $E(f)$ and rectifier distortion in $S(f)$ are greatly reduced as the frequencies for the carrier and twice the carrier are both above the phase-locking roll-off. (D) The strength of modulation coding was evaluated as the sum of the power near the first three harmonics of $F_m$ (gray boxes in panels A-C) for $S(f)$ (blue squares) and $E(f)$ (red circles). $VS_{pp}$ was also quantified to CF-tones for each fiber (black dashed line, right y-axis). (E) Rectifier distortion (RD) analysis was limited to the second harmonic of the carrier (brown boxes in panels A-C). RD was quantified as the sum of power in 10-Hz bands around twice the carrier frequency ($2 \times CF$) and the adjacent sidebands ($2 \times CF \pm F_m$). RD for $E(f)$ is not shown because $E(f)$ was virtually free from RD. (F) Raw and adjusted *sumcor* peak-heights across CFs. *sumcors* were adjusted by band-pass filtering them in the three $F_m$-related bands. Large differences between the two metrics at low frequencies indicate that the raw *sumcor* peak-heights are confounded by rectifier distortion at these frequencies. (G) Relation between raw and adjusted *sumcor* peak-heights with $F_m$-related power (from panel D) in $S(f)$. Good correspondence between $F_m$-related power in $S(f)$ and adjusted *sumcor* peak-height supports the use of spectrally specific envelope analyses. (H) The difference between raw and adjusted *sumcor* peak-heights was largely accounted for by RD power. However, this difference was always greater than zero, suggesting broadband metrics can also be biased because of noise related to neural stochasticity.

A second way envelope information in the neural response can be quantified is by computing the envelope of the difference PSTH, $d(t)$. This envelope, $e(t)$, can be estimated as the magnitude of the analytic signal, $a(t)$, of the difference PSTH

$$e(t) = \frac{|a(t)|}{\sqrt{2}}, \tag{6}$$

where $a(t) = d(t) + j\mathcal{H}\{d(t)\}$, and $\mathcal{H}\{\cdot\}$ is the Hilbert transform operator. The factor $\sqrt{2}$ normalizes for the power difference after applying the Hilbert transform. $d(t)$ is substantially less affected by rectifier distortion [23], and thus, so is $e(t)$. The use of $e(t)$ parallels the procedure followed by many computational models that extract envelopes from the output of cochlear filter banks [48, 54, 55].

The TFS component can also be estimated in two ways: (1) $d(t)$, and (2) cosine of the Hilbert phase of $d(t)$. The difference PSTH has been traditionally called the TFS response because it is the polarity-sensitive component. *difcor* and derived metrics relate to $d(t)$ as the *difcor* is related to the autocorrelation function of $d(t)$ (S4 Appendix). However, $d(t)$ does not represent the response to only the carrier (phase) since it also contains envelope information in $e(t)$. We propose a novel representation of the TFS response component, $\phi(t)$, estimated as the cosine phase of the analytic signal

$$\phi(t) = \sqrt{2} \times rms[d(t)] \times cos[\angle a(t)], \tag{7}$$

where normalization by $\sqrt{2} \times rms[d(t)]$ is used to match the power in $\phi(t)$ with the power in $d(t)$ since $cos[\angle a(t)]$ is a constant-rms ($rms = 1/\sqrt{2}$) signal.

**Relative merits of sum and Hilbert-envelope PSTHs in representing spike-train envelope responses.** The relative merits of the two envelope PSTHs, $s(t)$ and $e(t)$, were evaluated based on simulated spike-train data generated using a computational model of AN responses [56]. The model includes both cochlear-tuning and hair-cell transduction nonlinearities in the auditory system. Modulation spectra for sinusoidally amplitude-modulated (SAM) tones were estimated for $s(t)$ and $e(t)$ [denoted by $S(f)$ and $E(f)$, respectively] for individual-fiber responses (Fig 5A–5C). $d(t)$ was band-pass filtered near CF (200-Hz bandwidth, 2nd order filter) before applying the Hilbert transform to minimize the spectral energy in $d(t)$ that was not stimulus related. The two envelopes were evaluated based on their representations of the modulator and rectifier distortion. Rectifier distortions are expected to occur at even multiples of the carrier and nearby sidebands (i.e., $2nF_c$, $2nF_c - F_m$, and $2nF_c + F_m$ for integers n, Fig 5A). It is desirable for an envelope metric to consistently represent envelope coding across CFs and be less affected by rectifier-distortion artifacts. Modulation coding for the simulated responses was quantified as the power in 10-Hz bands centered at the first three harmonics of $F_m$ (i.e., 15 to 25 Hz, 35 to 45 Hz, and 55 to 65 Hz) for both $s(t)$ and $e(t)$ (Fig 5D). The need to include multiple harmonics of $F_m$ arises because the response during a stimulus cycle departs from sinusoidal shape due to the saturating nonlinearity associated with inner-hair-cell transduction (S2 Fig). While $F_m$-related power was nearly constant across CF for $s(t)$, it was nearly constant for $e(t)$ only up to 1.2 kHz, after which it rolled off. This roll-off for $e(t)$ is not surprising since $e(t)$ relies on phase-locking near the carrier and the sidebands, as confirmed by the strong correspondence between tonal phase-locking at the carrier frequency and $F_m$-related power in $e(t)$ (Fig 5D).

The analysis of rectifier distortion was limited to only the distortion components near the second harmonic of the carrier (i.e., $2F_c$, $2F_c - F_m$, and $2F_c + F_m$) since this harmonic is substantially stronger than higher harmonics (e.g., Fig 5A). Rectifier distortion was quantified as the sum of power in 10-Hz bands centered at the three distortion frequency components.

Because $e(t)$ was estimated from spectrally specific $d(t)$, which was band-limited to 200 Hz near the carrier frequency, $e(t)$ was virtually free from rectifier distortion. In contrast, $s(t)$ was substantially affected by rectifier distortion for simulated fibers with CFs below $\sim$2 kHz (Fig 5E). Rectifier distortion in $S(f)$ dropped for fibers with CF above $\sim$0.8 kHz because phase locking at distortion frequencies (i.e., twice the carrier frequencies) was attenuated by the roll-off in tonal phase locking. For example, the simulated AN fiber in Fig 5B (CF = 1.7 kHz) maintained comparable $F_m$-related power for both envelopes, but rectifier distortion for $s(t)$ was substantially diminished because the distortion frequency (3.4 kHz) is well above the phase-locking roll-off. These results indicate that $s(t)$ is substantially affected by rectifier distortion (at twice the stimulus frequency) when the neuron responds to stimulus energy that is below half the phase-locking cutoff.

Next, these spectral power metrics were compared with the correlogram-based metric, *sumcor* peak-height (Fig 5F–5H). The *sumcor* peak-height metric is defined as the maximum value of the normalized time-domain *sumcor* function [10]. Prior to estimating the peak-height, the *sumcor* is sometimes adjusted by adding an inverted triangular window to compensate for its triangular shape [13]. Here, *sumcors* were compensated by subtracting a triangular window from it so that the baseline *sumcor* is a flat function with a value of 0 (instead of 1) in the absence of ENV coding. In S5 Appendix, we show that the *sumcor* peak-height is a broadband metric and it is related to the total power in $s(t)$, including rectifier distortions. When the *sumcor* is used to analyze responses of low-frequency AN fibers to broadband noise stimuli, the *sumcor*-spectrum, and thus, the *sumcor* peak-height, are confounded by rectifier distortions [13]. Similar to $S(f)$ for low-frequency SAM tones (Fig 5A), these distortions show up at $2 \times$ CF in the *sumcor*-spectrum, whereas the *difcor*-spectrum has energy only near CF [13]. Heinz and colleagues addressed these distortions by low-pass filtering the *sumcor* below CF to remove the effects of rectifier distortion at $2 \times$ CF. Here, we generalize this issue by comparing the *sumcor* and spectrally specific ENV metrics for narrowband SAM-tone stimuli to demonstrate the limitations of any broadband ENV metric. *sumcors* were adjusted by band-limiting them to 10-Hz bands near the first three harmonics of $F_m$. As expected, the difference between the raw and adjusted *sumcor* peak-heights was large at low CFs (Fig 5F), where rectifier distortion corrupts the broadband *sumcor* peak-height metric. At high CFs (above 1.5 kHz), the difference between raw and adjusted *sumcor* peak-heights was small but nonzero. These differences correspond to power in $S(f)$ at frequencies other than the modulation-related bands and reflect the artifacts of neural stochasticity due to finite number of stimulus trials. As power is always nonnegative, including power at frequencies unrelated to the target frequencies adds bias and variance to any broadband metric. The adjusted *sumcor* peak-height, unlike the raw *sumcor* peak-height, showed good agreement with spectrally specific $F_m$-related power in $S(f)$ (Fig 5G).

Overall, these results support the use of spectrally specific analyses to quantify ENV coding in order to minimize artifacts due to rectifier distortion as well as the effects of neural stochasticity. Of the two candidate *apPSTHs* to quantify response envelope, $e(t)$ had the benefit of minimizing rectifier distortion. However, $e(t)$'s reliance on carrier-related phase locking limits the use of $e(t)$ as a unifying ENV metric across the whole range of CFs. Instead, spectrally specific $s(t)$ is more attractive because of its robustness in representing the response envelope across CFs (Fig 5D).

**Relative merits of difference and Hilbert-phase PSTHs in representing spike-train TFS responses.** In order to evaluate the relative merits of $d(t)$ and $\phi(t)$ in representing the neural TFS response, the same set of simulated AN spike-train responses were used as in Fig 5. Although the stimulus has power at the carrier ($F_c$) and sidebands ($F_c \pm F_m$; 6 dB lower), only the carrier representation should be considered towards quantifying the TFS response because
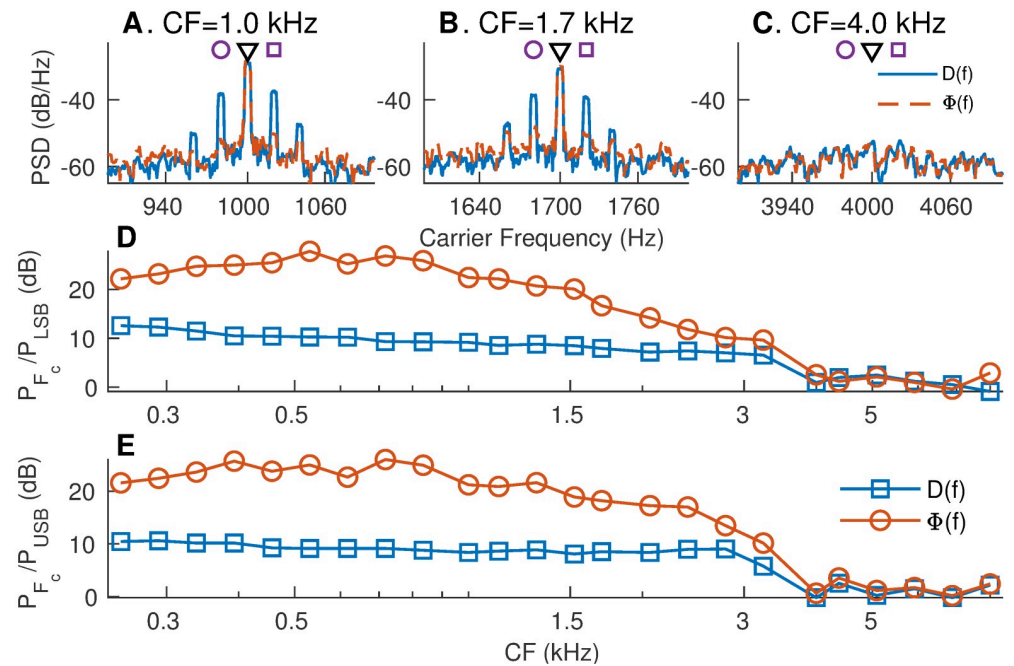
**Fig 6. Compared to the *d(t)*, the *apPSTH* $\phi(t)$ provides a better TFS representation.** (A-C) Spectra of $d(t)$ and $\phi(t)$ for the same three simulated AN fiber responses for which ENV spectra were shown in Fig 5. $D(f)$ has substantial power at CF (black triangle), as well as at lower (purple circle) and upper (purple square) sidebands. $\Phi(f)$, the spectrum of $\phi(t)$, shows maximum power concentration at *CF* (carrier frequency), with greatly reduced sidebands. (D) Ratio of power at CF (carrier, black triangle in panels A-C) to power at lower sideband (LSB, $F_c - F_m$, purple circles in panels A-C). (E) Ratio of power at CF (carrier) to power at upper sideband (USB, $F_c + F_m$, purple squares in panels A-C). $\phi(t)$ highlights the carrier and not the sidebands, and thus, compared to $d(t)$, $\phi(t)$ is a better representation of the true TFS response.

the energy at the sidebands arises due to the modulation of the carrier by the modulator (ENV). As the carrier has energy at a single frequency ($F_c$) for a SAM tone, the desirable TFS response should have maximum energy concentrated at the carrier frequency and not the sidebands. Therefore, the merits of $d(t)$ and $\phi(t)$ were evaluated based on how well they capture the carrier and suppress the sidebands (Fig 6).

As mentioned previously, $d(t)$ was band-limited to a 200-Hz bandwidth near the carrier frequency before estimating $\phi(t)$. $D(f)$ at low CFs contained substantial energy at both the carrier and the sidebands (Fig 6A and 6B). This indicates that $d(t)$ represents the complete neural coding of the SAM tone (both the envelope and the carrier) and not just the carrier. Furthermore, $D(f)$ has additional sidebands ($F_c \pm 2F_m$) around the carrier frequency. These sidebands arise as a result of the saturating nonlinearity associated with inner-hair-cell transduction (S2 Fig), and thus, should not be considered towards TFS response. In contrast, $\Phi(f)$, the spectrum of $\phi(t)$ had most of its power concentrated at the carrier frequency, with substantially less power in the sidebands (Fig 6A and 6B). These results were consistent across a wide range of CFs and for both sidebands (Fig 6D and 6E). Overall, these results show that $\phi(t)$ is a better PSTH compared to $d(t)$ in quantifying the response TFS since $\phi(t)$ emphasizes power at the carrier frequency and not at the sidebands.

In the following, we apply *apPSTH*-based analyses on spike-train data recorded from chinchilla AN fibers in response to speech and speech-like stimuli. In these examples, we particularly focus on certain ENV features, such as pitch coding for vowels and response onset for consonants, and TFS features, such as formant coding for vowels.
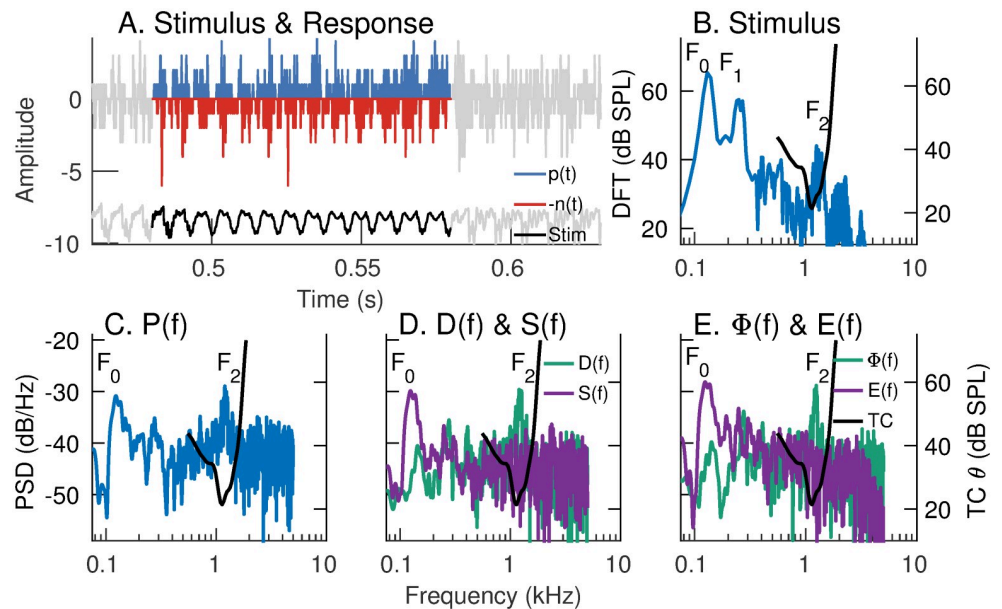
**Fig 7. Spectral-domain application of various *apPSTHs* to spike trains recorded in response to natural speech.**
Example of spectral analyses of spike trains recorded from an AN fiber (CF = 1.1 kHz, SR = 64 spikes/s) in response to a vowel snippet of a speech stimulus ($s_3$). (A) Time-domain representation of $p(t)$, $n(t)$, and the stimulus (*Stim*). $n(t)$ is reflected across the x-axis for display. Signals outside the analysis window are shown in gray. PSTH bin width = 0.1 ms. Number of stimulus repetitions per polarity = 50. Stimulus intensity = 65 dB SPL. (B) Stimulus spectrum (blue, left y-axis). In panels B-E, the frequency-threshold tuning curve (TC $\theta$, black) of the neuron is plotted on the right y-axis. (C) $P(f)$, which shows comparable energy at $F_0$ (130 Hz) and $F_2$ (1.2 kHz). (D) $D(f)$ and $S(f)$. (E) $\Phi(f)$ and $E(f)$. Both $S(f)$ and $E(f)$ show peaks near $F_0$. Similarly, both $D(f)$ and $\Phi(f)$ show good $F_2$ representations, although $D(f)$ is confounded by the strong $F_0$-related modulation in $e(t)$ as $d(t) = e(t) \times \phi(t)$. The significant representation of $F_0$ in this near-$F_2$ AN fiber response to a natural vowel is inconsistent with the synchrony-capture phenomenon for synthetic stationary vowels.

https://doi.org/10.1371/journal.pcbi.1008155.g007

**Neural characterization of ENV and TFS using *apPSTHs* for a natural speech segment.** Most previous studies have used the period histogram to study speech coding in the spectral domain [12, 57]. The period histogram is limited to stationary periodic stimuli, which were employed in those studies. In contrast, the use of *apPSTHs* facilitates the spectral analysis of neural responses to natural speech stimuli, which need not be stationary. Fig 7 shows the response spectra obtained using various *apPSTHs* [$p(t)$, $s(t)$, $d(t)$, and $\phi(t)$] for a low-frequency AN fiber in response to a natural speech segment (see S3 Fig for similar analyses for synthesized speech demonstrating the well-known "synchrony-capture" phenomenon [12, 57]). In this example, the response of a low-frequency AN fiber to a 100-ms vowel segment of the $s_3$ natural speech sentence was considered. The CF (1.1 kHz) of this neuron is close to the second formant ($F_2$) of this segment (Fig 7B). $P(f)$ shows peaks corresponding to $F_2$ (~1.2 kHz) and $F_0$ (~130 Hz, Fig 7C). Similar to S3 Fig, both $D(f)$ and $\Phi(f)$ show substantial energy near the formant closest to the neuron's CF. In contrast to S3 Fig, $S(f)$ [and $E(f)$] shows substantial energy near the fundamental frequency (inconsistent with synchrony capture). A detailed discussion of this discrepancy is beyond the scope of the present report, except to say that this lack of synchrony capture for natural vowels is a consistent finding that will be reported in a future study. The presence of substantial energy near $F_0$ in $E(f)$ indicates that $d(t)$ is confounded by pitch-related modulation in $e(t)$. This is because, mathematically, $D(f)$ is the convolution of the true TFS spectrum [$\Phi(f)$] and the Hilbert-envelope spectrum [$E(f)$]. Overall, these results
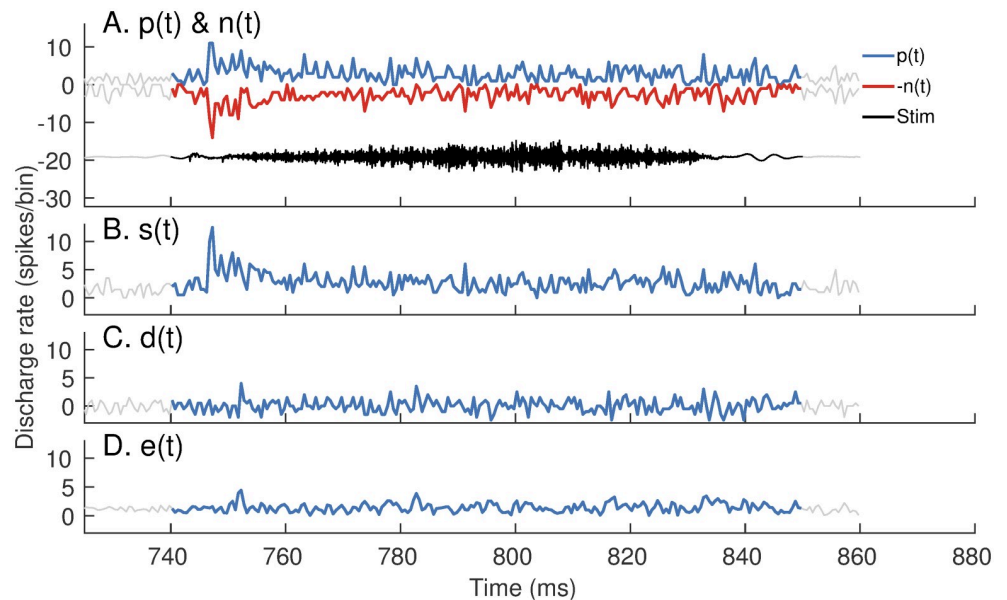
**Fig 8. $p(t)$, $n(t)$, and $s(t)$ have robust representations of the onset response, whereas $e(t)$ and $d(t)$ do not.** Response of a high-frequency fiber (CF = 5.8 kHz, SR = 70 spikes/s) to a fricative portion (/*s*/) of the speech stimulus, $s_3$. Stimulus intensity = 65 dB SPL. (*A*) Stimulus (black, labeled *Stim*), $p(t)$ (blue) and $n(t)$ (red, reflected across the x-axis). PSTH bin width = 0.5 ms. Number of stimulus repetitions per polarity = 50. (B) The sum envelope, $s(t)$ (C) The difference PSTH, $d(t)$, and (D) the Hilbert-envelope PSTH, $e(t)$. Since the onset envelope is a polarity-tolerant response, all PSTHs capture the response onset except for $d(t)$ and $e(t)$.

https://doi.org/10.1371/journal.pcbi.1008155.g008

demonstrate the application of various *apPSTHs* to study the neural representation of natural nonstationary speech stimuli in the spectral domain.

**Onset envelope is well represented in the sum PSTH but not in the Hilbert-envelope PSTH.**   In addition to analyzing spectral features, *apPSTHs* can also be used to analyze temporal features in the neural response. An example temporal feature is the onset envelope, which has been shown to be important for neural coding of consonants [22, 58], in particular fricatives [59]. A diminished onset envelope in the peripheral representation of consonants is hypothesized to be a contributing factor for perceptual deficits experienced by hearing-impaired listeners [60], and thus is important to quantify. Fig 8 shows example onset responses for a high-frequency AN fiber (CF = 5.8 kHz, SR = 70 spikes/s) for a fricative (/*s*/) portion of the speech stimulus $s_3$. The onset is well captured in single-polarity PSTHs [$p(t)$ and $n(t)$, Fig 8A] and in the sum envelope [$s(t)$, Fig 8B]. Since the onset is a polarity-tolerant feature, it is greatly reduced by subtracting the PSTHs to opposite polarities. As a result, response onset is poorly captured in $d(t)$ (Fig 8C) and its Hilbert envelope, $e(t)$ (Fig 8D).

Overall, these examples show that *apPSTHs* can be used to study various spectral and temporal features in neural responses for natural stimuli in the ENV/TFS dichotomy. These *apPSTHs* are summarized in Table 1 (and illustrated in S1 Fig).

## Quantifying ENV and TFS using *apPSTHs* for nonstationary signals

In the discussion so far, we have argued for using spectrally specific metrics to analyze neural responses to stationary stimuli. Another example where spectral specificity is needed is in evaluating the neural coding of nonstationary speech features (e.g., formant transitions). Speech is a nonstationary signal and conveys substantial information in its dynamic spectral trajectories

(e.g., Fig 1A). A number of studies have investigated the robustness of the neural representation of dynamic spectral trajectories using frequency glides and frequency-modulated tones as the stimulus [61–64]. These studies have usually employed a spectrogram analysis. While a spectrogram is effective for analyzing responses to nonstationary signals with unknown parameters, it does not explicitly incorporate information about the stimulus, which is often designed by the experimenter. Since the spectrogram relies on a narrow moving temporal window, it offers poor spectral resolution due to the time-frequency uncertainty principle. The same limitation applies to wavelet transforms that rely on segmenting the signal into shorter windows, even though window length varies across frequency. Instead of using these windowing-based analyses, frequency demodulation and filtering can be used together to estimate power along a spectrotemporal trajectory more accurately as described below. While this demodulation-based method has been described previously for other signals [65], we apply this method to natural speech and extend this approach to construct a new spectrally compact time-frequency representation called the *harmonicgram*. These spectrally specific analyses will facilitate more sensitive metrics to investigate the coding differences between nonstationary features in natural speech and extensively studied stationary features in synthetic speech.

**Frequency-demodulation-based spectrotemporal filtering.** First, we describe the spectrotemporal filtering technique using an example stimulus with dynamic spectral components (Fig 9). The 2-second-long stimulus consists of three spectrotemporal trajectories: (1) a stationary tone at 1.4 kHz, (2) a stationary tone at 2 kHz, and (3) a dynamic linear chirp that moves from 400 to 800 Hz over the stimulus duration. We are interested in estimating the
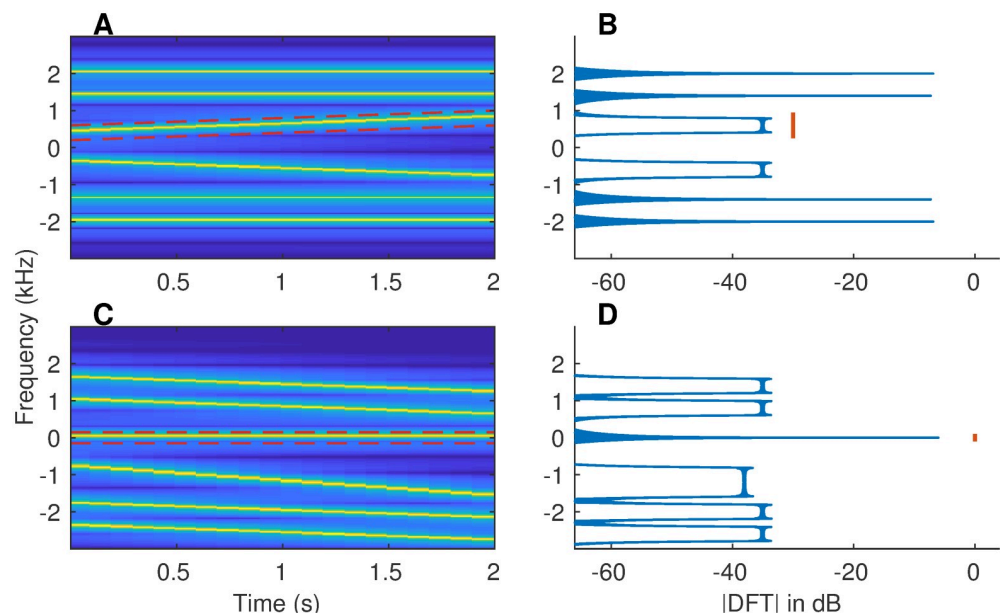


**Fig 9. More accurate estimates of power along a spectrotemporal trajectory can be obtained using frequency demodulation.** (A) Spectrogram of a synthesized example signal that mimics a single speech-formant transition. The 2-s signal consists of two stationary tones (1.4 and 2 kHz) and a linear frequency sweep (400 to 800 Hz). Red dashed lines outline the spectrotemporal trajectory along which we want to compute the power. Both positive and negative frequencies are shown for completeness. (B) Fourier-magnitude spectrum of the original signal. Energy related to the target spectrotemporal trajectory is spread over a wide frequency range (400 to 800 Hz, red line). (C) Spectrogram of the frequency-demodulated signal, where the target trajectory was used for demodulation (i.e., shifted down to 0 Hz). (D) Magnitude-DFT of the frequency-demodulated signal. The desired trajectory is now centered at 0 Hz, with its (spectral) energy spread limited only by the signal duration (i.e., equal to the inverse of signal duration), and hence, is much narrower.

power of the nonstationary component, the linear chirp. In order to estimate the power of this chirp, conventional spectrograms will employ one of the following two approaches. First, one can use a long window (e.g., 2 seconds) and compute power over the 400-Hz bandwidth from 400 to 800 Hz. In the second approach, one can use moving windows that are shorter in duration (e.g., 50 ms) and compute power with a resolution of 30 Hz (20-Hz imposed by inverse of the window duration and 10-Hz imposed by change in chirp frequency over 50 ms). As an alternative to these conventional approaches, one can demodulate the spectral trajectory of the linear chirp so that the chirp is demodulated to near 0 Hz (Fig 9C and 9D, see Materials and Methods). Then, a low-pass filter with 0.5-Hz bandwidth (as determined by the reciprocal of the 2-s stimulus duration) can be employed to estimate the time-varying power along the chirp trajectory. This time-varying power is estimated at the stimulus sampling rate, similar to the temporal sampling of the output of a band-pass filter applied on stationary signals. While the same temporal sampling can be achieved using the spectrogram by sliding the window by one sample and estimating the chirp-related power for each window, it will be computationally much more expensive compared to the frequency-demodulation-based approach. Furthermore, the spectral resolution of 0.5 Hz is the same as that for a stationary signal, which demonstrates a 60-fold improvement compared to the 50-ms window-based spectrogram approach.

**The *harmonicgram* for synthesized nonstationary speech.** As shown in Fig 9, combined use of frequency demodulation and low-pass filtering can provide an alternative to the spectrogram for analyzing signals with time-varying frequency components. Such an approach can also be used to study coding of dynamic stimuli that have harmonic spectrum with time-varying $F_0$, such as music and voiced speech. At any given time, a stimulus with a harmonic spectrum has substantial energy only at multiples of the fundamental frequency, $F_0$, which itself can vary with time [i.e., $F_0(t)$]. We take advantage of this spectral sparsity to introduce a new compact representation, the *harmonicgram*. Consider the $k$-th harmonic of $F_0(t)$; power along this trajectory [$kF_0(t)$] can be estimated using the frequency-demodulation-based spectrotemporal filtering technique. One could estimate the time-varying power along all integer multiples ($k$) of $F_0(t)$. This combined representation of the time-varying power across all harmonics of $F_0$ is the *harmonicgram* (see Materials and Methods). This name derives from the fact that this representation uses harmonic number instead of frequency (or spectrum) as in the conventional spectrogram.

Fig 10 shows harmonicgrams derived from *apPSTHs* in response to the nonstationary synthesized vowel, $s_2$. The first two formants are represented by their harmonic numbers, $F_1(t)/F_0(t)$ and $F_2(t)/F_0(t)$, which are known a priori in this case. Two harmonicgrams were constructed using responses from two AN fiber pools: (1) AN fibers that had a low CF (CF < 1 kHz), and (2) AN fibers that had a medium CF (1 kHz < CF < 2.5 kHz). Previous neurophysiological studies have shown that AN fibers with CF near and slightly above a formant strongly synchronize to that formant, especially at moderate to high intensities [12, 57]. Therefore, the low-CF pool was expected to capture $F_1$, which changed from 630 Hz to 570 Hz. Similarly, the medium-CF pool was expected to capture $F_2$, which changed from 1200 Hz to 1500 Hz. The harmonicgram for each pool was constructed by using the average Hilbert-phase PSTH, $\phi(t)$, of all AN fibers in the pool. The harmonicgram is shown from 38 ms to 188 ms to optimize the dynamic range to visually highlight the formant transitions by ignoring the onset response. The dominant component in the neural response for $F_1$ was expected at the harmonic number closest to $F_1/F_0$. For this stimulus, $F_1/F_0$ started at a value of 6.3 (630/100) and reached 4.75 (570/120) at 188 ms crossing 5.5 at 88.5 ms (Fig 10A). This transition of $F_1/F_0$ was faithfully represented in the harmonicgram where the dominant response switched from the 6th to the 5th harmonic near 90 ms. Similarly, $F_2/F_0$ started at 12, consistent with the dominant response at the 12th harmonic before 100 ms (Fig 10B). Towards the end of the stimulus, $F_2/F_0$ reached
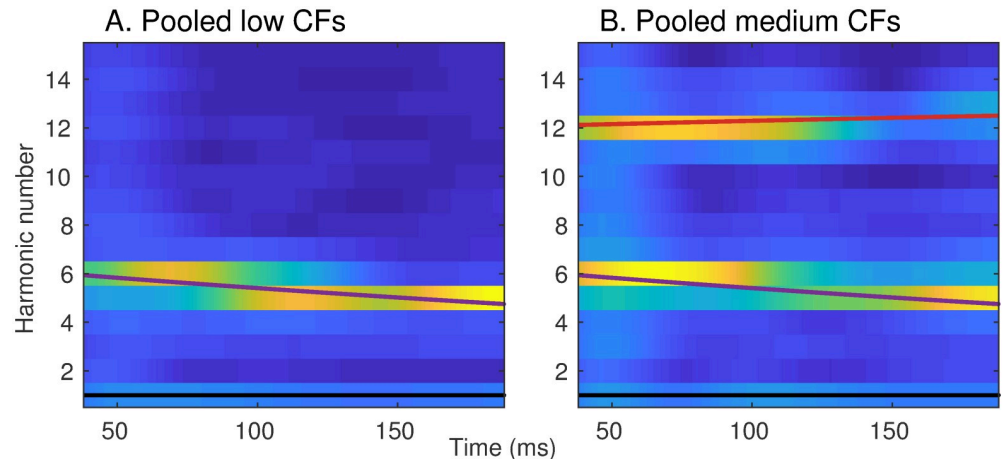
**Fig 10. The harmonicgram can be used to visualize formant tracking in synthesized nonstationary speech.** Neural harmonicgrams for fibers with a CF below 1 kHz (A, N = 16) and for fibers with a CF between 1 and 2.5 kHz (B, N = 29) in response to the dynamic vowel, $s_2$. Stimulus intensity = 65 dB SPL. The formant frequencies mimic formant trajectories of a natural vowel [21]. A 20-Hz bandwidth was employed to low-pass filter the demodulated signal for each harmonic. The harmonicgram for each AN-fiber pool was constructed by averaging the Hilbert-phase PSTHs of all AN fibers within the pool. PSTH bin width = 50 $\mu$s. Data are from one chinchilla. The black, purple, and red lines represent the fundamental frequency ($F_0/F_0$), the first formant ($F_1/F_0$) and the second formant ($F_2/F_0$) contours, respectively. The time-varying formant frequencies were normalized by the time-varying $F_0$ to convert the spectrotemporal representation into a harmonicgram.

12.5, which is consistent with the near-equal power in the 12th and the 13th harmonic in the harmonicgram. In contrast to findings from previous studies, the harmonicgram for the medium-CF pool indicates that these fibers respond to both the first and second formants [57, 66]. Such a complex response with components corresponding to multiple formants is likely due to the steep slope of the vowel spectrum (S4 Fig).

**The harmonicgram for natural speech.** The harmonicgram analysis is not limited to synthesized vowels, but can also be applied to natural speech (Fig 11). These harmonicgrams were constructed for the natural speech stimulus, $s_3$, using average $\phi(t)$ for the same low-CF and medium-CF AN fiber pools used in Fig 10. Here, we consider a 500-ms segment of the stimulus, which contains multiple phonemes. Qualitatively, similar to Fig 10, these harmonicgrams capture formant contours across phonemes. The harmonicgram for the low-CF pool emphasizes the $F_1$ contour, whereas the harmonicgram for the medium-CF pool primarily emphasizes the $F_2$ contour, and to a lesser extent, the $F_1$ contour. Compared to the spectrogram, the harmonicgram representation is more compact and spectrally specific. Furthermore, from a neural-coding perspective, quantifying how individual harmonics of $F_0$ are represented in the response is more appealing than the spectrogram since response energy is concentrated only at these $F_0$ harmonics.

The harmonicgram not only provides a compact representation for nonstationary signals with harmonic spectra, it can also be used to quantify coding strength of time-varying features, such as formants for speech (Fig 11E and 11F). In these examples, the strength of formant coding at each time point, $t$, was quantified as the sum of power in the three harmonics closest to the $F_0$-normalized formant frequency at that time [e.g., $F_1(t)/F_0(t)$]. As expected, power for the harmonics near the first formant was substantially greater than for the second formant for the low-CF pool (Fig 11E). For the medium-CF pool, $F_2$ representation was robust over the whole stimulus duration, although $F_1$ representation was largely comparable (Fig 11F). These
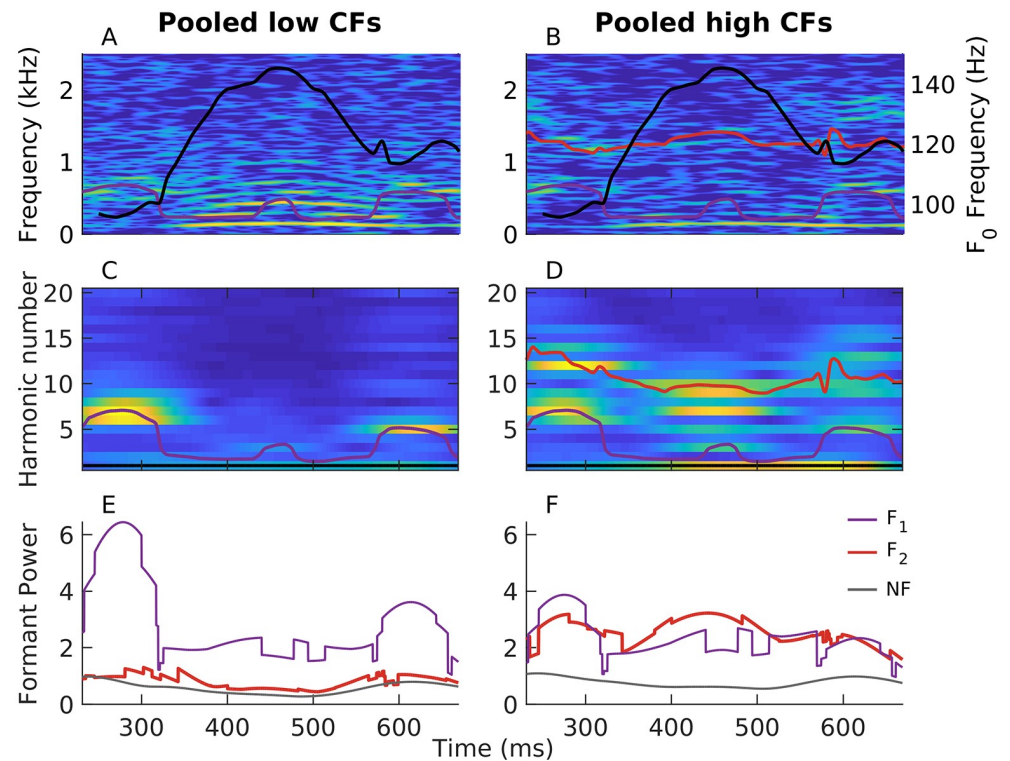
**Fig 11. The harmonicgram can be used to quantify the coding of time-varying stimulus features at superior spectrotemporal resolution compared to the spectrogram.** Harmonicgrams were constructed using $\phi(t)$ for the same two AN-fiber pools described in Fig 10. PSTH bin width = 50 $\mu$s. A 9-Hz bandwidth was employed to low-pass filter the demodulated signal for each harmonic. The data were collected from one chinchilla in response to the speech stimulus, $s_3$. Stimulus intensity = 65 dB SPL. A 500-ms segment corresponding to the voiced phrase "amle" was considered. (A, B) Spectrograms constructed from the average $\phi(t)$ for the low-CF pool (A) and from the medium-CF pool (B). (C, D) Average harmonicgrams for the same set of fibers as in A and B, respectively. Warm (cool) colors represent regions of high (low) power. The first-formant contour ($F_1$ in A and B, $F_1/F_0$ in C and D) is highlighted in purple. The second-formant contour ($F_2$ in A and B, $F_2/F_0$ in C and D) is highlighted in red. Trajectories of the fundamental frequency (black in A and B, right y-axis) and the formants were obtained using Praat [67]. (E, F) Harmonicgram power near the first formant (purple) and the second formant (red) for the low-CF pool (E) and the medium-CF pool (F). Harmonicgram power for each formant at any given time ($t$) was computed by summing the power in the three closest $F_0$ harmonics adjacent to the normalized formant contour [e.g., $F_1(t)/F_0(t)$] at that time. The noise floor (NF) for power was estimated as the sum of power for the 29th, 30th, and 31st harmonics of $F_0$ because the frequencies corresponding to these harmonics were well above the CFs of both fiber pools. These time-varying harmonicgram power metrics are spectrally specific to $F_0$ harmonics and are computed with high temporal sampling rate (same as the original signal). This spectrotemporal resolution is much better than the spectrotemporal resolution that can be obtained using spectrograms.

https://doi.org/10.1371/journal.pcbi.1008155.g011

examples demonstrate novel analyses using the *apPSTH*-based harmonicgram to quantify time-varying stimulus features in single-unit neural responses at high spectrotemporal resolution, which is not possible with conventional windowing-based approaches.

**The harmonicgram can also be used to analyze FFRs in response to natural speech.** As mentioned earlier, a major benefit of using *apPSTHs* to analyze spike trains is that the same analyses can also be applied to evoked far-field potentials. In Fig 12, the harmonicgram analysis was applied to the difference FFR recorded in response to the same speech sentence ($s_3$) that was used in Fig 11. In fact, these FFR data and spike-train data used in Fig 11 were collected from the same chinchilla. The difference FFR was computed as the difference between FFRs to opposite polarities of the stimulus. The spectrogram and harmonicgram can also be constructed using the Hilbert-phase FFR to highlight the TFS component of the response (S5
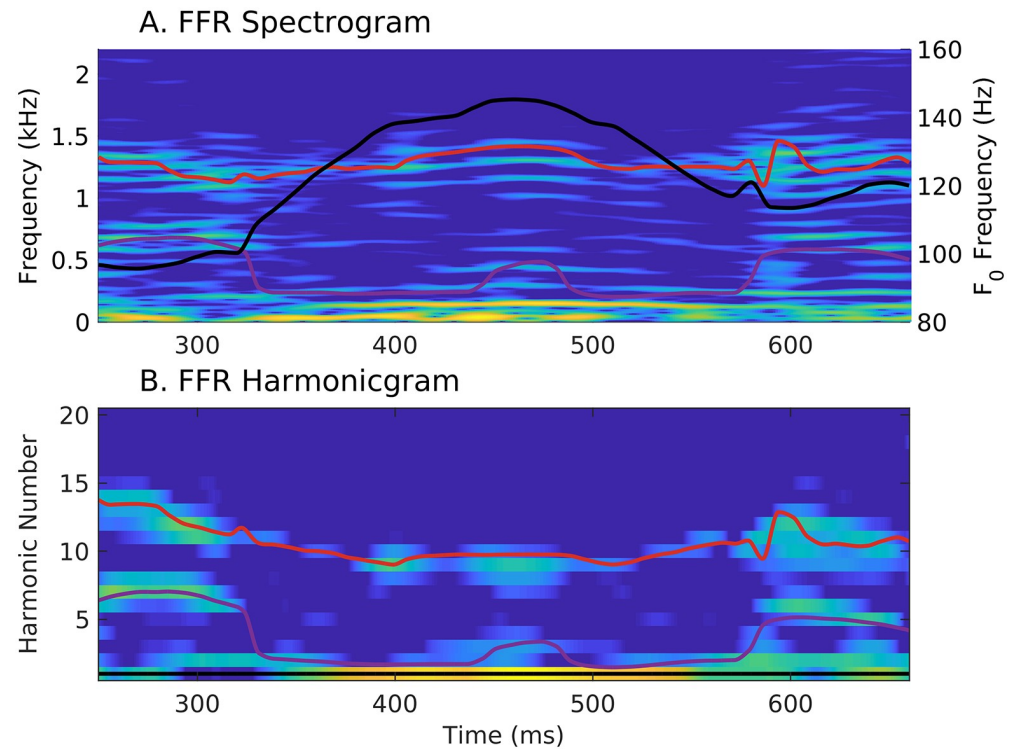
**Fig 12. The harmonicgram of the FFR to natural speech shows robust dynamic tracking of formant trajectories, similar to the AN-fiber harmonicgram.** Comparison of the spectrogram (A) and the harmonicgram (B) for the FFR recorded in response to the same stimulus, $s_3$ that was used to analyze *apPSTHs* in Fig 11. Stimulus intensity = 65 dB SPL. Lines and colormap are the same as in Fig 11. These plots were constructed using the difference FFR, which reflects the neural coding of both stimulus TFS and ENV. To highlight the coding of stimulus TFS, Hilbert-phase $[\phi(t)]$ FFR can be used instead of the difference FFR (S5 Fig). The FFR harmonicgram (A) is strikingly similar to the AN-fiber harmonicgrams in Fig 11C and 11D in that the representations of the first two formants are robust. The FFR data here and spike-train data used in Fig 11 were obtained from the same animal.

https://doi.org/10.1371/journal.pcbi.1008155.g012

Fig). Unlike the *apPSTHs* for AN fibers, the FFR cannot be used to construct two sets of harmonicgrams corresponding to different populations of neurons because the FFR lacks tonotopic specificity. Nevertheless, this FFR-harmonicgram is strikingly similar to the medium-CF pool harmonicgram in Fig 11D. The dynamic representations of the first two formants are robust in both the representations. In fact, the FFR representations seem more robust in formant tracking compared to PSTH-derived representations, qualitatively, especially for the harmonicgram. A more uniform sample of neurons contribute to evoked responses compared to the AN fiber sample corresponding to Fig 11, which could be a factor for the robustness of the FFR representations. Overall, these results reinforce the idea that using *apPSTHs* to analyze spike trains offers the same spectrally specific analyses that can be applied to evoked far-field potentials, e.g., the FFR, thus allowing a unifying framework to study temporal coding for both stationary and nonstationary signals in the auditory system.

## Discussion

### Use of *apPSTHs* underlies a unifying framework to study temporal coding in the auditory system

A better understanding of the neural correlates of perception requires the integration of electrophysiological, psychophysical, and neurophysiological analyses in the same framework.

Although extensive literature exists in both electrophysiology and neurophysiology on the neural correlates of perception, the analyses employed in these studies have diverged. This disconnect is largely because the forms of the neural data are different (i.e., continuous-valued waveforms versus point-process spike trains). The present report provides a unifying framework for analyzing spike trains using *apPSTHs*, which offers numerous benefits over previous neurophysiological analyses. Specifically, the use of *apPSTHs* incorporates many of the previous ad-hoc approaches, such as VS and correlograms (Eqs 3 to 5). In fact, correlograms and metrics derived from them can be estimated using *apPSTHs* in a computationally efficient way. The *apPSTHs* essentially convert the naturally rectified neurophysiological point-process data into a continuous-valued signal, which allows advanced signal processing tools designed for continuous-valued signals to be applied to spike-train data. For example, *apPSTHs* can be used to derive spectrally specific TFS components [e.g., $\phi(t)$, Fig 6], multitaper spectra (Fig 3), modulation-domain representations (Fig 4), and harmonicgrams (Figs 10 and 11). *apPSTHs* can also be directly compared to evoked far-field responses for both stationary and nonstationary stimuli (e.g., Figs 11 and 12).

## Temporal coding metrics should be spectrally specific

The various analyses explored here advocate for spectral specificity of temporal coding metrics. The need for spectrally specific analyses arises for two reasons: (1) neural data is finite and stochastic, and (2) spike-train data are rectified. Neural stochasticity exacerbates spectral-estimate variance at all frequencies; therefore, time-domain (equivalently broadband) metrics will be noisier compared to narrowband metrics. Similarly, the rectified nature of spike-train data introduces harmonic distortions in the response spectrum, which can corrupt broadband metrics (e.g., TFS distortion at two times the carrier frequency corrupting estimates of ENV coding, Fig 5A and 5B).

These issues requiring spectral specificity are not unique to the *apPSTH* analyses but also apply to classic metrics, e.g., correlograms. For example, the broadband correlation index (CI) metric is appropriate to analyze responses of neurons with high CFs, but the CI metric is confounded by rectifier distortions for neurons with low CFs [11, 13]. Studies have previously tried to avoid these distortions in the *sumcor* by restricting the response bandwidth to below the CF because, for a given filter, the envelope bandwidth cannot be greater than the filter bandwidth [13, 68].

Here, we have extended and generalized the analysis of these issues using narrowband stimuli. In particular, when a neuron responds to low-frequency stimulus energy that is below half the phase-locking cutoff, responses that contain any polarity-tolerant component [e.g., $p(t)$, $n(t)$, $s(t)$, SAC, and *sumcor*] will be confounded by rectifier distortion of the polarity-sensitive component (Fig 5E). Any broadband metric of temporal coding should exclude these distortions at twice the carrier frequency. Beyond avoiding rectifier distortion, limiting the bandwidth of a metric to only the desired bands will lead to more precise analyses by minimizing the effects of neural stochasticity (Fig 5H). For example, envelope coding metrics for SAM-tone stimuli should consider the spectrum power only at $F_m$ and its harmonics [69], rather than the simple approach of always low-pass filtering at CF [13].

Similar to envelope-based metrics, metrics that quantify TFS coding should also be spectrally specific to the carrier frequency. Previous metrics of TFS coding, such as $d(t)$ and *difcor*, are not specific to the carrier frequency but rather include modulation sidebands as well as additional sidebands due to transduction nonlinearities (Fig 6). In contrast, $\phi(t)$ introduced here emphasizes the carrier and suppresses the sidebands (Fig 6). Thus, the spectrally specific

$\phi(t)$ is a better TFS response, which relates to the zero-crossing signal used in the signal processing literature [70–72].

## Spectral-estimation benefits of using *apPSTHs*

Neurophysiological studies have usually favored the DFT to estimate the response spectrum. For example, the DFT has been applied to the period histogram [12, 57], the single-polarity PSTH [73, 74], the difference PSTH [23], and correlograms [10]. Since spike-train data are stochastic and usually sparse and finite, there is great scope for spectral estimates, including the DFT spectrum, to suffer from bias and variance issues. The multitaper approach optimally uses the available data to minimize the bias and variance of the spectral estimate [4, 39, 75]. The multitaper approach can be used with both *apPSTHs* and correlograms, but using *apPSTHs* offers additional variance improvement up to a factor of 2 (Fig 3). This improvement is because twice as many tapers (both odd and even) can be used with an *apPSTH* compared to a correlogram, which is an even sequence and limits analyses to only using even tapers. Additional benefits may be achievable by combining the Lomb-Scargle approach, which is well-suited for estimating the spectrum of unevenly sampled data (e.g., spike trains), with *apPSTHs* in the multitaper framework [76].

## Benefits of spectrotemporal filtering

Analysis of neural responses to nonstationary signals has been traditionally carried out using windowing-based approaches, such as the spectrogram. Shorter windows help with tracking rapid temporal structures, but they offer poorer spectral resolution. On the other hand, larger windows allow better spectral resolution at the cost of smearing rapid dynamic features. As an alternative to windowing-based approaches, spectrotemporal filtering can improve the spectral resolution of analyses by taking advantage of stimulus parameters that are known a priori (Fig 9). This approach is particularly efficient to analyze spectrally sparse signals (i.e., signals with instantaneous line spectra, such as voiced speech). In particular, the spectral resolution is substantially improved compared to the spectrogram. In addition, while the same temporal sampling can be obtained using the spectrogram, it will be much more computationally expensive compared to the spectrotemporal filtering approach, as discussed in the following example.

The benefits of spectrotemporal filtering extend to other spectrally sparse signals, like harmonic complexes. A priori knowledge of the fundamental frequency can be used to construct the harmonicgram, which takes advantage of power concentration at harmonics of $F_0$. This approach contrasts with the spectrogram, which computes power at all frequencies uniformly. The harmonicgram can be used to analyze both kinematic synthesized vowels (Fig 10) as well as natural speech (Fig 11). The harmonicgram is particularly useful in quantifying dominant harmonics at high temporal sampling and is thus applicable to nonstationary signals. The harmonicgram can also be applied to evoked far-field potentials (e.g., the FFR in Fig 12). While alternatives exist to analyze spike-train data in response to time-varying stimuli [77], the present spectrotemporal technique is simpler and can be directly applied to both spike-train data and far-field responses. Overall, these results support the idea that using *apPSTHs* to analyze spike trains provides a unifying framework to study temporal coding in the auditory system across modalities. Furthermore, this framework facilitates the study of dynamic-stimulus coding by the nonlinear and time-varying auditory system.

### *apPSTHs* allow animal models of sensorineural hearing loss to be linked to psychophysical speech-intelligibility models

Speech-intelligibility models not only improve our understanding of perceptually relevant speech features, but they can also be used to optimize hearing-aid and cochlear-implant strategies. However, existing SI models work well for normal-hearing listeners but have not been widely extended for hearing-impaired listeners. This gap is largely because of the fact that most SI models are based on signal-processing algorithms in the acoustic domain, where individual differences in the physiological effects of various forms of sensorineural hearing loss on speech coding are difficult to evaluate. This gap can be addressed by extending acoustic SI models to the neural spike-train domain. In particular, spike-train data obtained from preclinical animal models of sensorineural hearing loss can be used to explore the neural correlates of perceptual deficits faced by hearing-impaired listeners [78]. These insights will be crucial for developing accurate SI models for hearing-impaired listeners.

*apPSTHs* offer a straightforward means to study various speech features in the neural spike-train domain. As *apPSTHs* are in the same discrete-time continuous-valued form as acoustic signals, acoustic SI models can be directly translated to the neural domain. Many successful SI models are based on the representation of the temporal envelope [16, 48], although the role of TFS remains a matter of controversy [79]. In fact, recent studies suggest that the peripheral representation of TFS can shape central envelope representations, and thereby alter speech perception outcomes [80, 81]. *apPSTHs* can be used to derive modulation-domain representations so that envelope-based SI models can be evaluated in the neural domain (Fig 4). Similarly, the Hilbert-phase PSTH, $\phi(t)$, can be used to evaluate the neural representation of TFS features. These TFS results will be particularly insightful for cochlear-implant stimulation strategies that rely on the zero-crossing component of the stimulus, which closely relates to $\phi(t)$ [82, 83].

### Translational benefits of animal models

A key motivation of this paper was to develop a framework so that insights and findings from animal models can ultimately improve our understanding of how the human auditory system processes real-life sounds, like speech. Experiments involving human subjects are typically limited to far-field responses, such as compound action potentials, frequency-following responses, and auditory brainstem responses. However, these evoked responses include contributions from multiple sources such as the cochlear microphonic, electrical interferences, and responses from several neural substrates [34, 37]; these contributions are not clearly understood. The *apPSTH*-based framework offers a straightforward way to study these contributions by comparing anatomically specific spike-train responses with clinically viable noninvasive responses.

This framework is also beneficial to develop and validate noninvasive metrics using animal models and apply these metrics to humans. For example, we demonstrated the applicability of the new spectrally compact harmonicgram approach on both spike-train data and FFR data recorded from chinchillas to evaluate speech coding. This harmonicgram analysis can also be applied to FFR data recorded from humans to study natural speech coding in both normal and impaired auditory systems. Similarly, the representation of other important response features, such as the onset and adaptation, can also be linked between invasive and noninvasive data using preclinical animal models of different forms of SNHL. Overall, these insights will be informative for estimating the anatomical and physiological states of humans using noninvasive measures, and how these states relate to individual differences in speech perception that currently challenge audiological rehabilitation.

### Limitations

**Biological feasibility.**   The analyses proposed here aim to rigorously quantify the dichotomous ENV/TFS information in the neural response and bridge the definitions between the audio and neural spike-train domains. Methods discussed here may not all be biologically feasible. For example, the brain does not have access to both polarities of the stimulus. Thus, the PSTHs that require two polarities to be estimated, e.g., $s(t)$, $d(t)$, and $\phi(t)$, may not have an "internal representation" in the brain. This limitation also applies to correlogram metrics based on *sumcor* and *difcor*, which require two polarities of the stimulus. Thus, the use of the single-polarity PSTH [$p(t)$] to derive the central "internal representations" is more appropriate from a biological feasibility perspective (e.g., Fig 4). However, these various ENV/TFS components allow a thorough characterization of the processing of spectrotemporally complex signals by the nonlinear auditory system and can guide the development of more accurate speech-intelligibility models and help improve signal processing strategies for hearing-impaired listeners.

**Alternating-polarity stimuli.**   Use of two polarities may not be sufficient to separate out all components underlying neural responses when more than two components contribute to neural responses at a given frequency. In particular, it may be intractable to separate out rectifier distortion when the bandwidths of ENV and TFS responses overlap. For example, consider the response of a broadly tuned AN fiber to a vowel, which has a fundamental frequency of $F_0$. The energy at $2F_0$ in $S(f)$ may reflect one or more of the following sources: (1) rectifier distortion to carrier energy at $F_0$, (2) beating between (carrier) harmonics that are separated by $2F_0$, and (3) effects of transduction nonlinearities on the beating between (carrier) harmonics that are separated by $F_0$. In these special cases, additional stimulus phase variations can be used to separate out these components [84, 85].

**The harmonicgram.**   A key drawback of applying the harmonicgram to natural speech is the requirement of knowing the $F_0$ trajectory. $F_0$ estimation is a difficult problem, especially in degraded speech. Thus, the harmonicgram could be inaccurate unless the $F_0$ trajectory is known, or at least the original stimulus is known so that $F_0$ can be estimated. A second confound is the unknown stimulus-to-response latency for different systems. Latencies for different neurons vary with their CF, stimulus frequency, and stimulus intensity. Thus, even if the acoustic spectrotemporal trajectory is known precisely, errors may accumulate if latencies are not properly accounted for. This issue will likely be minor for spectrotemporal trajectories with slow dynamics. For stimuli with faster dynamics, latency confounds can be easily minimized by estimating stimulus-to-response latency by cross-correlation and using a larger cutoff frequency for low-pass filtering.

## Materials and methods

### Ethics statement

All procedures followed NIH-issued guidelines and were approved by the Purdue Animal Care and Use Committee (Protocol No: 1111000123).

### Experimental procedures

Spike trains were recorded from single AN fibers of anesthetized chinchillas using standard procedures in our laboratory [68, 86]. Anesthesia was induced with xylazine (2 to 3 mg/kg, subcutaneous) and ketamine (30 to 40 mg/kg, intraperitoneal), and supplemented with sodium pentobarbital ($\sim$7.5 mg/kg/hour, intraperitoneal). FFRs were recorded using subdermal electrodes in a vertical montage (mastoid to vertex with common ground near the nose)

under the same ketamine/xylazine anesthesia induction protocol described above using standard procedures in our laboratory [87]. Spike times were stored with 10-$\mu$s resolution. FFRs were stored with 48-kHz sampling rate. Stimulus presentation and data acquisition were controlled by custom MATLAB-based (The MathWorks, Natick, MA) software that interfaced with hardware modules from Tucker-Davis Technologies (TDT, Alachua, FL) and National Instruments (NI, Austin, TX).

## Speech stimuli

The following four stimuli were used in these experiments. ($s_1$) Stationary vowel, $\wedge$ (as in c<u>u</u>p; S1 Audio): $F_0$ was 100 Hz. The first three formants were placed at $F_1 = 600$, $F_2 = 1200$, and $F_3 = 2500$ Hz. The vowel was 188 ms in duration. ($s_2$) Nonstationary vowel, $\wedge$ (S2 Audio): $F_0$ increased linearly from 100 to 120 Hz over its 188-ms duration. The first two formants moved as well ($F_1$: 630 $\rightarrow$ 570 Hz; $F_2$: 1200 $\rightarrow$ 1500 Hz; see S4 Fig). $F_3$ was fixed at 2500 Hz. The formant frequencies for both $s_1$ and $s_2$ were chosen based on natural formant contours of the vowel $\wedge$ in American English [21, 88]. $s_1$ and $s_2$ were synthesized using a MATLAB instantiation of the Klatt synthesizer (courtesy of Dr. Michael Kiefte, Dalhousie University, Canada). ($s_3$) A naturally uttered Danish sentence (list #1, sentence #3 in the CLUE Danish speech intelligibility test, [89]). ($s_4$) A naturally uttered English sentence (Sentence #2, List #1 in the Harvard Corpus, [90]). All speech and speech-like stimuli were played at an overall intensity of 60 to 65 dB SPL.

## Power along a spectro-temporal trajectory

Consider a known frequency trajectory, $f_{traj}(t)$, along which we need to estimate power in a signal, $x(t)$. The phase trajectory, $\Phi traj(t)$, can be computed as

$$\Phi_{traj}(t) = \int_0^t f_{traj}(\tau)d\tau. \tag{8}$$

For discrete-time signals, the phase trajectory can be estimated as

$$\Phi_{traj}[n] = \frac{1}{f_s}\sum_{m=1}^{n} f_{traj}[m]. \tag{9}$$

The phase trajectory can be demodulated from $x(t)$ by multiplying a complex exponential with phase = $-\Phi_{traj}(t)$ [65]

$$x_{demod}(t) = x(t) \ e^{-j2\pi\Phi_{traj}(t)}. \tag{10}$$

The power along $f_{traj}(t)$ in $x(t)$ can be estimated as the power in $x_{demod}(t)$ within the spectral-resolution bandwidth (W) near 0 Hz in the spectral estimate, $P_{x_{demod}}(f)$.

$$P_{traj} = 2\int_{-W/2}^{W/2} P_{x_{demod}}(f)df. \tag{11}$$

The scaling factor 2 is required because the integral in Eq 11 only represents the original positive-frequency band of the real signal, $x(t)$; the equal amount of power within the original negative-frequency band, which is shifted further away from 0 Hz by $\Phi_{traj}(t)$, should also be included (see Fig 9).

## The harmonicgram

Consider a harmonic complex, $x(t)$, with a time-varying (instantaneous) fundamental frequency, $F_0(t)$. For a well-behaved and smooth $F_0(t)$, energy in $x(t)$ will be concentrated at multiples of the instantaneous fundamental frequency, i.e., $kF_0(t)$. Thus, $x(t)$ can be represented by the energy distributed across the harmonics of the fundamental. The time-varying power along the $k$-th harmonic of $F_0(t)$ can be estimated by first demodulating $x(t)$ with the $kF_0(t)$ trajectory using Eq 10, and then using an appropriate low-pass filter to limit energy near 0 Hz (say within $\pm W/2$). We define the *harmonicgram* as the matrix of time-varying power along all harmonics of the fundamental frequency. Thus, the harmonicgram is

$$harmonicgram(k, t) = \mathcal{LPF}_{[-W/2, W/2]}\{x(t)\ e^{-j2\pi kF_0(t)}\}. \tag{12}$$

## Supporting information

**S1 Text. Classic metrics for quantifying temporal coding in the auditory system.**
(PDF)

**S1 Audio. Stimulus 1 ($s_1$).** Stationary synthesized vowel, $\wedge$.
(WAV)

**S2 Audio. Stimulus 2 ($s_2$).** Nonstationary synthesized vowel, $\wedge$.
(WAV)

**S1 Appendix. Vector strength metric definitions.**
(PDF)

**S2 Appendix. Relation between the *vector strength* metric and the *difference* PSTH.**
(PDF)

**S3 Appendix. Relation between *shuffled correlograms* and *apPSTHs*.**
(PDF)

**S4 Appendix. Relation between *difcor/sumcor* and *difference/sum* PSTHs.**
(PDF)

**S5 Appendix. Relation between *shuffled-correlogram* peak-height and *apPSTHs*.**
(PDF)

**S1 Table. Glossary of terms and definitions.**
(PDF)

**S2 Table. Parameters for the AN model.**
(PDF)

**S1 Fig. Graphical illustration of *apPSTHs* in Table 1.**
(PDF)

**S2 Fig. Nonlinear inner-hair-cell transduction function introduces additional sidebands in the spectrum for a SAM tone.**
(PDF)

**S3 Fig. Neural characterization of ENV and TFS using *apPSTHs* for a synthesized stationary vowel.**
(PDF)

**S4 Fig. DFT-magnitude for the nonstationary vowel, $s_2$.**
(PDF)

**S5 Fig. FFR harmonicgram can be constructed using the Hilbert-phase response.**
(PDF)

## Acknowledgments

## Author Contributions

**Conceptualization:** Satyabrata Parida, Hari Bharadwaj, Michael G. Heinz.

**Data curation:** Satyabrata Parida, Michael G. Heinz.

**Formal analysis:** Satyabrata Parida.

**Funding acquisition:** Michael G. Heinz.

**Investigation:** Satyabrata Parida.

**Methodology:** Satyabrata Parida, Hari Bharadwaj, Michael G. Heinz.

**Software:** Satyabrata Parida.

**Visualization:** Satyabrata Parida.

**Writing – original draft:** Satyabrata Parida.

**Writing – review & editing:** Satyabrata Parida, Hari Bharadwaj, Michael G. Heinz.

## References

1. Tremblay KL, Billings CJ, Friesen LM, Souza PE. Neural Representation of Amplified Speech Sounds. Ear and Hearing. 2006; 27(2):93–103. https://doi.org/10.1097/01.aud.0000202288.21315.bd

2. Clinard CG, Tremblay KL, Krishnan AR. Aging alters the perception and physiological representation of frequency: Evidence from human frequency-following response recordings. Hearing Research. 2010; 264(1):48–55.

3. Kraus N, Anderson S, White-Schwoch T. The Frequency-Following Response: A Window into Human Communication. In: Kraus N, Anderson S, White-Schwoch T, Fay RR, Popper AN, editors. The Frequency-Following Response: A Window into Human Communication. Springer Handbook of Auditory Research. Cham: Springer International Publishing; 2017. p. 1–15.

4. Thomson DJ. Spectrum estimation and harmonic analysis. Proceedings of the IEEE. 1982; 70 (9):1055–1096. https://doi.org/10.1109/PROC.1982.12433

5. Moore BC. Cochlear hearing loss: physiological, psychological and technical issues. John Wiley & Sons; 2007.

6. Scharenborg O. Reaching over the gap: A review of efforts to link human and automatic speech recognition research. Speech Communication. 2007; 49(5):336–347. https://doi.org/10.1016/j.specom.2007.01.009

7. Goldberg JM, Brown PB. Response of binaural neurons of dog superior olivary complex to dichotic tonal stimuli: some physiological mechanisms of sound localization. Journal of Neurophysiology. 1969; 32(4):613–636. https://doi.org/10.1152/jn.1969.32.4.613

8. Rees A, Palmer AR. Neuronal responses to amplitude-modulated and pure-tone stimuli in the guinea pig inferior colliculus, and their modification by broadband noise. The Journal of the Acoustical Society of America. 1989; 85(5):1978–1994. https://doi.org/10.1121/1.397851

9. Joris PX, Yin TCT. Responses to amplitude-modulated tones in the auditory nerve of the cat. The Journal of the Acoustical Society of America. 1992; 91(1):215–232. https://doi.org/10.1121/1.402757

**10.** Louage DHG, Heijden Mvd, Joris PX. Temporal Properties of Responses to Broadband Noise in the Auditory Nerve. Journal of Neurophysiology. 2004; 91(5):2051–2065. https://doi.org/10.1152/jn.00816.2003

**11.** Joris PX, Louage DH, Cardoen L, van der Heijden M. Correlation Index: A new metric to quantify temporal coding. Hearing Research. 2006; 216–217:19–30.

**12.** Young ED, Sachs MB. Representation of steady-state vowels in the temporal aspects of the discharge patterns of populations of auditory-nerve fibers. The Journal of the Acoustical Society of America. 1979; 66(5):1381–1403. https://doi.org/10.1121/1.383532

**13.** Heinz MG, Swaminathan J. Quantifying Envelope and Fine-Structure Coding in Auditory Nerve Responses to Chimaeric Speech. Journal of the Association for Research in Otolaryngology. 2009; 10 (3):407–423. https://doi.org/10.1007/s10162-009-0169-8

**14.** Colburn HS, Carney LH, Heinz MG. Quantifying the Information in Auditory-Nerve Responses for Level Discrimination. Journal of the Association for Research in Otolaryngology. 2003; 4(3):294–311. https://doi.org/10.1007/s10162-002-1090-6

**15.** Paliwal KK, Alsteris L. Usefulness of Phase Spectrum in Human Speech Perception. Eighth European Conference on Speech Communication and Technology. 2003; p. 4.

**16.** Relaño-Iborra H, May T, Zaar J, Scheidiger C, Dau T. Predicting speech intelligibility based on a correlation metric in the envelope power spectrum domain. The Journal of the Acoustical Society of America. 2016; 140(4):2670–2679. https://doi.org/10.1121/1.4964505

**17.** Heil P, Peterson AJ. Basic response properties of auditory nerve fibers: a review. Cell Tissue Res. 2015; 361(1):129–158. https://doi.org/10.1007/s00441-015-2177-9

**18.** Sayles M, Heinz MG. Afferent Coding and Efferent Control in the Normal and Impaired Cochlea. In: Understanding the Cochlea. Springer Handbook of Auditory Research. Springer, Cham; 2017. p. 215–252.

**19.** Nearey TM, Assmann PF. Modeling the role of inherent spectral change in vowel identification. The Journal of the Acoustical Society of America. 1986; 80(5):1297–1308. https://doi.org/10.1121/1.394433

**20.** Delgutte B. Auditory neural processing of speech. The handbook of phonetic sciences. 1997; p. 507–538.

**21.** Hillenbrand JM, Nearey TM. Identification of resynthesized /hVd/ utterances: Effects of formant contour. The Journal of the Acoustical Society of America. 1999; 105(6):3509–3523. https://doi.org/10.1121/1.424676

**22.** Delgutte B. Representation of speech-like sounds in the discharge patterns of auditory-nerve fibers. The Journal of the Acoustical Society of America. 1980; 68(3):843–857. https://doi.org/10.1121/1.384824

**23.** Sinex DG, Geisler CD. Responses of auditory-nerve fibers to consonant–vowel syllables. The Journal of the Acoustical Society of America. 1983; 73(2):602–615. https://doi.org/10.1121/1.389007

**24.** Cariani PA, Delgutte B. Neural correlates of the pitch of complex tones. I. Pitch and pitch salience. Journal of Neurophysiology. 1996; 76(3):1698–1716. https://doi.org/10.1152/jn.1996.76.3.1698

**25.** Sayles M, Winter IM. Reverberation Challenges the Temporal Representation of the Pitch of Complex Sounds. Neuron. 2008; 58(5):789–801. https://doi.org/10.1016/j.neuron.2008.03.029

**26.** Shannon RV, Zeng FG, Kamath V, Wygonski J, Ekelid M. Speech Recognition with Primarily Temporal Cues. Science. 1995; 270(5234):303–304. https://doi.org/10.1126/science.270.5234.303

**27.** Smith ZM, Delgutte B, Oxenham AJ. Chimaeric sounds reveal dichotomies in auditory perception. Nature. 2002; 416(6876):87. https://doi.org/10.1038/416087a

**28.** Liberman MC. Auditory-nerve response from cats raised in a low-noise chamber. The Journal of the Acoustical Society of America. 1978; 63(2):442–455. https://doi.org/10.1121/1.381736

**29.** Goblick TJ, Pfeiffer RR. Time-Domain Measurements of Cochlear Nonlinearities Using Combination Click Stimuli. The Journal of the Acoustical Society of America. 1969; 46(4B):924–938. https://doi.org/10.1121/1.1911812

**30.** Kemp DT. Stimulated acoustic emissions from within the human auditory system. The Journal of the Acoustical Society of America. 1978; 64(5):1386–1391. https://doi.org/10.1121/1.382104

**31.** Aiken SJ, Picton TW. Envelope and spectral frequency-following responses to vowel sounds. Hearing Research. 2008; 245(1):35–47.

**32.** Shinn-Cunningham B, Ruggles DR, Bharadwaj H. How Early Aging and Environment Interact in Everyday Listening: From Brainstem to Behavior Through Modeling. In: Moore BCJ, Patterson RD, Winter IM, Carlyon RP, Gockel HE, editors. Basic Aspects of Hearing. Advances in Experimental Medicine and Biology. Springer New York; 2013. p. 501–510.

**33.** Ananthakrishnan S, Krishnan A, Bartlett E. Human Frequency Following Response: Neural Representation of Envelope and Temporal Fine Structure in Listeners with Normal Hearing and Sensorineural Hearing Loss. Ear and Hearing. 2016; 37(2):e91–e103. https://doi.org/10.1097/AUD.0000000000000247

**34.** King A, Hopkins K, Plack CJ. Differential Group Delay of the Frequency Following Response Measured Vertically and Horizontally. Journal of the Association for Research in Otolaryngology. 2016; 17(2):133–143. https://doi.org/10.1007/s10162-016-0556-x

**35.** Joris PX. Interaural Time Sensitivity Dominated by Cochlea-Induced Envelope Patterns. Journal of Neuroscience. 2003; 23(15):6345–6350. https://doi.org/10.1523/JNEUROSCI.23-15-06345.2003

**36.** Lichtenhan JT, Cooper NP, Guinan JJ. A new auditory threshold estimation technique for low frequencies: Proof of concept. Ear and Hearing. 2013; 34(1):42–51. https://doi.org/10.1097/AUD.0b013e31825f9bd3

**37.** Verschooten E, Joris PX. Estimation of Neural Phase Locking from Stimulus-Evoked Potentials. Journal of the Association for Research in Otolaryngology. 2014; 15(5):767–787. https://doi.org/10.1007/s10162-014-0465-9

**38.** Parida S, Heinz MG. Noninvasive Measures of Distorted Tonotopic Speech Coding Following Noise-Induced Hearing Loss. Journal of the Association for Research in Otolaryngology. 2020;. https://doi.org/10.1007/s10162-020-00755-2 PMID: 33188506

**39.** Babadi B, Brown EN. A Review of Multitaper Spectral Analysis. IEEE Transactions on Biomedical Engineering. 2014; 61(5):1555–1564. https://doi.org/10.1109/TBME.2014.2311996

**40.** Oppenheim AV. Discrete-time signal processing. Pearson Education India; 1999.

**41.** Rangayyan RM. Biomedical signal analysis. vol. 33. John Wiley & Sons; 2015.

**42.** Kryter KD. Methods for the Calculation and Use of the Articulation Index. The Journal of the Acoustical Society of America. 1962; 34(11):1689–1697. https://doi.org/10.1121/1.1909094

**43.** Houtgast T, Steeneken HJM. The Modulation Transfer Function in Room Acoustics as a Predictor of Speech Intelligibility. Acta Acustica united with Acustica. 1973; 28(1):66–73.

**44.** Taal CH, Hendriks RC, Heusdens R, Jensen J. An Algorithm for Intelligibility Prediction of Time–Frequency Weighted Noisy Speech. IEEE Transactions on Audio, Speech, and Language Processing. 2011; 19(7):2125–2136. https://doi.org/10.1109/TASL.2011.2114881

**45.** Cooke M. A glimpsing model of speech perception in noise. The Journal of the Acoustical Society of America. 2006; 119(3):1562–1573. https://doi.org/10.1121/1.2166600

**46.** Heinz MG. Neural modelling to relate individual differences in physiological and perceptual responses with sensorineural hearing loss. Proceedings of the International Symposium on Auditory and Audiological Research. 2015; 5:137–148.

**47.** Rallapalli VH, Heinz MG. Neural Spike-Train Analyses of the Speech-Based Envelope Power Spectrum Model: Application to Predicting Individual Differences with Sensorineural Hearing Loss. Trends in Hearing. 2016; 20, 1–14.

**48.** Jørgensen S, Dau T. Predicting speech intelligibility based on the signal-to-noise envelope power ratio after modulation-frequency selective processing. The Journal of the Acoustical Society of America. 2011; 130(3):1475–1487. https://doi.org/10.1121/1.3621502

**49.** Scheidiger C, Carney LH, Dau T, Zaar J. Predicting Speech Intelligibility Based on Across-Frequency Contrast in Simulated Auditory-Nerve Fluctuations. Acta Acustica united with Acustica. 2018; 104 (5):914–917. https://doi.org/10.3813/AAA.919245

**50.** Delgutte B, Hammond BM, Cariani PA. Neural coding of the temporal envelope of speech: relation to modulation transfer functions. Psychophysical and physiological advances in hearing. 1998; p. 595–603.

**51.** Steeneken HJM, Houtgast T. A physical method for measuring speech-transmission quality. The Journal of the Acoustical Society of America. 1980; 67(1):318–326. https://doi.org/10.1121/1.384464

**52.** Elhilali M, Chi T, Shamma SA. A spectro-temporal modulation index (STMI) for assessment of speech intelligibility. Speech Communication. 2003; 41(2):331–348.

**53.** Jørgensen S, Ewert SD, Dau T. A multi-resolution envelope-power based model for speech intelligibility. The Journal of the Acoustical Society of America. 2013; 134(1):436–446. https://doi.org/10.1121/1.4807563

**54.** Dubbelboer F, Houtgast T. The concept of signal-to-noise ratio in the modulation domain and speech intelligibility. The Journal of the Acoustical Society of America. 2008; 124(6):3937–3946. https://doi.org/10.1121/1.3001713

**55.** Sadjadi SO, Hansen JHL. Hilbert envelope based features for robust speaker identification under reverberant mismatched conditions. In: 2011 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP); 2011. p. 5448–5451.

**56.** Bruce IC, Erfani Y, Zilany MSA. A phenomenological model of the synapse between the inner hair cell and auditory nerve: Implications of limited neurotransmitter release sites. Hearing Research. 2018; 360:40–54. https://doi.org/10.1016/j.heares.2017.12.016

**57.** Delgutte B, Kiang NYS. Speech coding in the auditory nerve: I. Vowel-like sounds. The Journal of the Acoustical Society of America. 1984; 75(3):866–878. https://doi.org/10.1121/1.390596

**58.** Heil P. Coding of temporal onset envelope in the auditory system. Speech Communication. 2003; 41 (1):123–134. https://doi.org/10.1016/S0167-6393(02)00099-7

**59.** Delgutte B, Kiang NYS. Speech coding in the auditory nerve: III. Voiceless fricative consonants. The Journal of the Acoustical Society of America. 1984; 75(3):887–896. https://doi.org/10.1121/1.390598

**60.** Allen JB, Li F. Speech perception and cochlear signal processing [Life Sciences]. IEEE Signal Processing Magazine. 2009; 26(4):73–77. https://doi.org/10.1109/MSP.2009.932564

**61.** Krishnan A, Parkinson J. Human Frequency-Following Response: Representation of Tonal Sweeps. Audiology and Neurotology. 2000; 5(6):312–321. https://doi.org/10.1159/000013897

**62.** Skoe E, Kraus N. Auditory brainstem response to complex sounds: a tutorial. Ear and Hearing. 2010; 31(3):302–324. https://doi.org/10.1097/AUD.0b013e3181cdb272

**63.** Clinard CG, Cotter CM. Neural representation of dynamic frequency is degraded in older adults. Hearing Research. 2015; 323:91–98. https://doi.org/10.1016/j.heares.2015.02.002

**64.** Billings CJ, Bologna WJ, Muralimanohar RK, Madsen BM, Molis MR. Frequency following responses to tone glides: Effects of frequency extent, direction, and electrode montage. Hearing Research. 2019; 375:25–33. https://doi.org/10.1016/j.heares.2019.01.012

**65.** Olhede S, Walden At. A generalized demodulation approach to time-frequency projections for multi-component signals. Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences. 2005; 461(2059):2159–2179. https://doi.org/10.1098/rspa.2005.1455

**66.** Miller RL, Schilling JR, Franck KR, Young ED. Effects of acoustic trauma on the representation of the vowel /ɛ/ in cat auditory nerve fibers. The Journal of the Acoustical Society of America. 1997; 101 (6):3602–3616. https://doi.org/10.1121/1.418321

**67.** Boersma P. Praat, a system for doing phonetics by computer. Glot Int. 2001; 5(9):341–345.

**68.** Kale S, Heinz MG. Envelope Coding in Auditory Nerve Fibers Following Noise-Induced Hearing Loss. Journal of the Association for Research in Otolaryngology. 2010; 11(4):657–673. https://doi.org/10.1007/s10162-010-0223-6

**69.** Vasilkov V, Verhulst S. Towards a differential diagnosis of cochlear synaptopathy and outer-hair-cell deficits in mixed sensorineural hearing loss pathologies. medRxiv. 2019. https://doi.org/10.1101/19008680

**70.** Voelcker HB. Toward a unified theory of modulation—Part II: Zero manipulation. Proceedings of the IEEE. 1966; 54(5):735–755. https://doi.org/10.1109/PROC.1966.4843

**71.** Logan BF. Information in the Zero Crossings of Bandpass Signals. Bell System Technical Journal. 1977; 56(4):487–510. https://doi.org/10.1002/j.1538-7305.1977.tb00522.x

**72.** Wiley R. Approximate FM Demodulation Using Zero Crossings. IEEE Transactions on Communications. 1981; 29(7):1061–1065. https://doi.org/10.1109/TCOM.1981.1095091

**73.** Miller MI, Sachs MB. Representation of stop consonants in the discharge patterns of auditory-nerve fibers. The Journal of the Acoustical Society of America. 1983; 74(2):502–517. https://doi.org/10.1121/1.389816

**74.** Carney LH, Geisler CD. A temporal analysis of auditory-nerve fiber responses to spoken stop consonant–vowel syllables. The Journal of the Acoustical Society of America. 1986; 79(6):1896–1914. https://doi.org/10.1121/1.393197

**75.** Percival DB, Walden AT. Spectral analysis for physical applications. Cambridge University Press; 1993.

**76.** Springford A, Eadie GM, Thomson DJ. Improving the Lomb–Scargle Periodogram with the Thomson Multitaper. The Astronomical Journal. 2020; 159(5):205. https://doi.org/10.3847/1538-3881/ab7fa1

**77.** Brown EN, Barbieri R, Ventura V, Kass RE, Frank LM. The Time-Rescaling Theorem and Its Application to Neural Spike Train Data Analysis. Neural Computation. 2002; 14(2):325–346. https://doi.org/10.1162/08997660252741149

**78.** Trevino M, Lobarinas E, Maulden AC, Heinz MG. The chinchilla animal model for hearing science and noise-induced hearing loss. The Journal of the Acoustical Society of America. 2019; 146(5): 3710–3732. https://doi.org/10.1121/1.5132950

**79.** Lorenzi C, Gilbert G, Carn H, Garnier S, Moore BCJ. Speech perception problems of the hearing impaired reflect inability to use temporal fine structure. PNAS. 2006; 103(49):18866–18869. https://doi.org/10.1073/pnas.0607364103

80. Ding N, Chatterjee M, Simon JZ. Robust cortical entrainment to the speech envelope relies on the spectro-temporal fine structure. NeuroImage. 2014; 88:41–46. https://doi.org/10.1016/j.neuroimage.2013.10.054

81. Viswanathan V, Bharadwaj HM, Shinn-Cunningham B, Heinz MG. Evaluating human neural envelope coding as the basis of speech intelligibility in noise. The Journal of the Acoustical Society of America. 2019; 145(3):1717–1717. https://doi.org/10.1121/1.5101298

82. Grayden DB, Burkitt AN, Kenny OP, Clarey JC, Paolini AG, Clark GM. A cochlear implant speech processing strategy based on an auditory model. In: Proceedings of the 2004 Intelligent Sensors, Sensor Networks and Information Processing Conference, 2004.; 2004. p. 491–496.

83. Chen F, Zhang YT. Zerocrossing-based nonuniform sampling to deliver low-frequency fine structure cue for cochlear implant. Digital Signal Processing. 2011; 21(3):427–432. https://doi.org/10.1016/j.dsp.2010.12.002

84. Billings SA, Zhang H. Analysing non-linear systems in the frequency domain–II. The phase response. Mechanical Systems and Signal Processing. 1994; 8(1):45–62. https://doi.org/10.1006/mssp.1994.1004

85. Lucchetti F, Deltenre P, Avan P, Giraudet F, Fan X, Nonclercq A. Generalization of the primary tone phase variation method: An exclusive way of isolating the frequency-following response components. The Journal of the Acoustical Society of America. 2018; 144(4):2400–2412. https://doi.org/10.1121/1.5063821

86. Henry KS, Sayles M, Hickox AE, Heinz MG. Divergent auditory-nerve encoding deficits between two common etiologies of sensorineural hearing loss. Journal of Neuroscience. 2019; 39:6879–6887. https://doi.org/10.1523/JNEUROSCI.0038-19.2019

87. Zhong Z, Henry KS, Heinz MG. Sensorineural hearing loss amplifies neural coding of envelope information in the central auditory system of chinchillas. Hearing Research. 2014; 309:55–62. https://doi.org/10.1016/j.heares.2013.11.006

88. Hillenbrand J, Getty LA, Clark MJ, Wheeler K. Acoustic characteristics of American English vowels. The Journal of the Acoustical Society of America. 1995; 97(5):3099–3111. https://doi.org/10.1121/1.411872

89. Nielsen JB, Dau T. Development of a Danish speech intelligibility test. International Journal of Audiology. 2009; 48(10):729–741. https://doi.org/10.1080/14992020903019312

90. Rothauser EH. IEEE recommended practice for speech quality measurements. IEEE Trans on Audio and Electroacoustics. 1969; 17:225–246. https://doi.org/10.1109/TAU.1969.1162058