

# emg2speech: synthesizing speech from electromyography using self-supervised speech models

Anonymous

## Abstract

We present a neuromuscular speech interface that translates electromyographic (EMG) signals recorded from orofacial muscles during speech articulation directly into audio. We find that self-supervised speech representations (SS) are strongly linearly related to the electrical power of muscle activity: a simple linear mapping predicts EMG power from SS with a correlation of  $r = 0.85$ . In addition, EMG power vectors associated with distinct articulatory gestures form structured, separable clusters. Together, these observations suggest that SS implicitly encode articulatory mechanisms, as reflected in EMG activity. Leveraging this structure, we map EMG signals into the SS space and synthesize speech, enabling end-to-end EMG-to-speech generation without explicit articulatory modeling or vocoder training. We demonstrate this system with a participant with amyotrophic lateral sclerosis (ALS), converting orofacial EMG recorded while she *silently* articulated speech into audio.

 [PROJECT PAGE](#).  [GITHUB](#)<sup>1</sup>.

## 1 Introduction

Neural and neuromuscular interfaces hold significant promise for augmenting human abilities to interact and communicate with the external world. Brain-computer interfaces (BCIs), such as the speech neuroprostheses described in [Wairagkar et al. \(2025\)](#); [Metzger et al. \(2023\)](#); [Willett et al. \(2023\)](#), have shown that individuals with conditions such as anarthria or amyotrophic lateral sclerosis can regain functional speech through invasive neural recordings. While these invasive approaches are well suited for individuals with severe paralysis or complete loss of articulatory control, their widespread deployment is limited by the need for

surgical implantation, high cost, and clinical risk. In contrast, we propose a non-invasive speech interface that leverages preserved articulatory muscle activity, enabling a broader range of individuals, including those with laryngectomy, dysarthria, or dysphonia, to regain functional speech without surgical intervention.

In this article, we present a method that leverages self-supervised speech (SS) models to convert electromyographic (EMG) signals collected during speech articulation directly into audio, without explicitly training a vocoder. Our key insight comes from the observation that speech features derived from SS models can be linearly mapped to the electrical power of muscle action potentials. Because EMG power patterns associated with different articulatory gestures form structured and separable clusters in feature space, these results suggest that SS models implicitly encode articulatory information, as reflected in EMG activity. This relationship also motivates EMG power as an effective intermediate representation for mapping muscle activity to speech features. We exploit this property to design a lightweight EMG-to-audio conversion model that leverages EMG power representations in conjunction with SS models.

## 2 Prior work

Converting non-speech signals into audio has been explored across several modalities, including lip-movements-to-speech ([Kim et al., 2021](#); [Prajwal et al., 2020](#)), motor-cortex neural signals-to-speech ([Wairagkar et al., 2025](#); [Metzger et al., 2023](#); [Littlejohn et al., 2025](#)), and EMG-to-speech ([Gaddy and Klein, 2020, 2021](#)). Most existing approaches in these domains ([Kim et al., 2021](#); [Prajwal et al., 2020](#); [Wairagkar et al., 2025](#); [Gaddy and Klein, 2020, 2021](#)) assume that the alignment between the input signals (e.g., video or neural activity) and the corresponding audio is known. In contrast,

<sup>1</sup>We will open-source all data and model checkpoints upon publication of the manuscript. Due to IRB restrictions, we are unable to release the dataset anonymously.

we address a more challenging scenario, similar to Metzger et al. (2023); Littlejohn et al. (2025), where the alignment between neural activity (in our case, EMG) and speech is *unknown*. This setting requires the model not only to learn the mapping between EMG activity and audio but also to infer the underlying alignment from an exponential search space.

Work in Littlejohn et al. (2025); Metzger et al. (2023) addresses this alignment-free setting by training an encoder that takes motor-cortex neural signals as input and maps them to discrete HUBERT units (Lakhotia et al., 2021), which are then passed to a pretrained vocoder (Tacotron (Wang et al., 2017)) following the pipeline in Lakhotia et al. (2021). We adopt a similar high-level pipeline for EMG-to-speech conversion. However, our approach explicitly leverages the *geometric structure* of EMG signals and their relationship to self-supervised (SS) speech representations to design an encoder grounded in articulatory mechanisms.

Despite this progress, prior work faces several practical limitations. For instance, Littlejohn et al. (2025); Metzger et al. (2023) use a small-vocabulary corpus containing only 1,024 words, and in Littlejohn et al. (2025) (where motor-cortex neural activity is converted into speech), each test sentence was presented to the model an average of 6.94 times during training. Moreover, these studies are not fully reproducible due to limited implementation details and the lack of public access to the data used in the experiments. Since non-speech-to-speech conversion typically involves multiple components in an end-to-end pipeline, opaque designs make it difficult to reproduce results and to compare methods fairly. These limitations motivate the need for richer public datasets and reproducible benchmarks for fair evaluation, both of which we provide in this work.

## 2.1 Our contributions

We make three primary contributions.

**First**, we open-source one of the largest high quality EMG-to-speech datasets to date, comprising a large-vocabulary corpus of approximately 9 hours of EMG speech data with over 6,800 unique words from a healthy participant, as well as a small-vocabulary corpus of approximately 1 hour of EMG speech data with roughly 300 unique words from a participant with ALS. To the best of our knowledge, these datasets constitute one of the largest and most comprehensive publicly available resources

for EMG-to-speech conversion.

**Second**, building on these datasets, we develop encoder architectures grounded in articulatory mechanisms that exhibit interpretability and operate effectively in low-data regimes, including settings with as little as 40 minutes of training data from an ALS participant. This is particularly important given the practical difficulty of collecting large-scale EMG datasets with current sensing technology. To support learning under limited supervision, we leverage massively pretrained self-supervised speech models and use their representations as a structured target space for EMG-to-speech mapping. We further establish, for the first time, a quantitative relationship between EMG signals and self-supervised speech representations, and exploit this structure to guide encoder design.

**Third**, we introduce phoneme-guided decoding for EMG-to-speech synthesis, demonstrating that incorporating phonetic structure improves the quality of generated audio.

Unlike prior EMG-to-speech benchmarks (Gaddy and Klein, 2020, 2021), our approach does not assume known temporal alignment between EMG and audio during training. This design is motivated by clinically relevant scenarios in which parallel EMG-audio pairs may be unavailable or unreliable, such as laryngectomy (absence of laryngeal voicing) or ALS (degraded acoustic recordings due to bulbar impairment). As a result, the model must learn without frame-level EMG-audio correspondence, which substantially increases the difficulty of the learning problem. Overall, our contributions address fundamental aspects of EMG-to-speech modeling and are simple to implement, yet work well with widely available off-the-shelf pretrained components.

Because our study derives audio from text transcripts rather than using time-aligned EMG-audio pairs, and targets an unaligned EMG-to-speech synthesis setting, there are no existing benchmarks that support direct one-to-one comparisons. Nevertheless, where possible, we compare against representative baselines from recent EMG interface literature, including spectrogram-based feature pipelines from EMG2QWERTY (Sivakumar et al., 2024), to contextualize performance. In appendix C, we additionally provide broader comparisons to prior brain-computer speech interfaces for context, although these are not intended as direct one-to-one comparisons.

### 3 Data

We use three datasets in this study: (i) a large, general-corpus vocabulary dataset from a healthy participant, denoted  $\text{DATA}_{\text{GENERAL}}$ ; (ii) a small, limited-corpus vocabulary dataset from a participant with ALS, denoted  $\text{DATA}_{\text{ALS}}$ ; and (iii) a dataset of discrete orofacial gestures underlying speech articulation collected from 12 healthy participants, denoted  $\text{DATA}_{\text{OROFACIAL GESTURES}}$ . Below, we describe each dataset in detail.

#### 3.1 $\text{DATA}_{\text{GENERAL}}$

A healthy participant naturally articulates English sentences while the corresponding EMG signals are recorded at 5000 Hz. We record EMG from 31 muscle sites on the neck, chin, jaw, cheek, and lips using monopolar surface electrodes (see figure 5 for electrode placement and appendix A for additional details).

We adapt the language corpus from Willett et al. (2023), who demonstrate a speech brain-computer interface by translating motor-cortex activity into speech. The dataset comprises an English corpus with approximately 6800 unique words and 9660 sentences. Sentences vary in length, and the participant articulates at a normal speaking rate, averaging 115 words per minute. We split the dataset into training, validation, and test sets containing 8500, 760, and 400 sentences, respectively. Sentences in the test set do not appear in the training or validation sets.

The start and end of each sentence are time-stamped using mouse clicks. When the participant is ready to begin, they click the mouse to display the sentence on the screen and mark the start time. After articulation is complete, they click again to mark the end time; this second click removes the sentence from the screen. This procedure allows the participant to articulate at their own pace.

#### 3.2 $\text{DATA}_{\text{ALS}}$

A participant diagnosed with amyotrophic lateral sclerosis (ALS) silently articulates English sentences (with overt articulatory movements but no audible output) while we record the corresponding EMG signals at 5000 Hz using the same electrode layout as in  $\text{DATA}_{\text{GENERAL}}$ .

We construct a small English corpus comprising approximately 300 unique words and 600 sentences. Sentences vary in length, and the participant articulates at her current comfortable speak-

ing rate, averaging 61 words per minute. We split the dataset into training, validation, and test sets containing 500, 40, and 60 sentences, respectively. Test sentences do not appear in the training or validation sets.

#### 3.3 $\text{DATA}_{\text{OROFACIAL GESTURES}}$

Twelve participants perform 13 distinct orofacial movements, with 10 repetitions per movement. The set of movements includes cheeks: puff out, cheeks: suck in, jaw: drop down, jaw: move backward, jaw: move forward, jaw: move left, jaw: move right, lips: pucker, lips: smile, lips: tuck (as if blotting), tongue: back of lower teeth, tongue: back of upper teeth, and tongue: roof of the mouth. These movements are selected to span a broad range of articulatory gestures involved in natural speech production, including lip rounding, jaw positioning, and tongue placement, which are essential for producing different phonemes.

This dataset is recorded using 22 electrodes at a sampling rate of 5000 Hz. Signals are recorded from approximately the same muscle sites as in  $\text{DATA}_{\text{GENERAL}}$ , except that electrodes on the right side of the neck are not used (middle diagram in figure 5). Each gesture is performed for a duration of 1.5 s.

## 4 Methods

### 4.1 Electromyography (EMG)

EMG signals are collected by a set of sensors  $\mathcal{V}$  and represented as functions of time  $t$ . A sequence of EMG signals  $E$  corresponding to articulated speech, associated with an audio signal  $A$  and phonemic content  $L$ , is represented as  $E = \{\mathbf{f}_v(t)\}_{v \in \mathcal{V}}$ . Here,  $\mathbf{f}_v(t)$  denotes the EMG signal captured at sensor node  $v$  as a function of time. The audio signal  $A$  encodes both phonemic (lexical) content and expressive aspects of speech such as volume, pitch, prosody, and intonation, while  $L$  represents only the phonemic content—a sequence of phonemes. For example, the phonemic content  $L$  of the word <FRIDAY> is denoted by the phoneme sequence <F-R-IY-D-AY>.

**EMG covariance matrices:** for an EMG signal  $E_{\mathcal{V} \times \tau}$  collected from  $\mathcal{V}$  sensor nodes over a duration of  $\tau$  samples, we construct a symmetric positive definite (SPD) covariance matrix  $\mathcal{E}_{\mathcal{V} \times \mathcal{V}} = \epsilon E E^\top$ , where  $\epsilon$  is a scaling factor. We denote the diagonal of  $\mathcal{E}$  as  $\mathbb{D}(\mathcal{E})$  and its lower triangular part

as  $\lfloor \mathcal{E} \rfloor$ . The vector  $\mathbb{D}(\mathcal{E})$  represents the muscle action potential power at each electrode  $v \in \mathcal{V}$  during the interval  $\tau$ , while the off-diagonal elements capture the pairwise cross-channel covariance, reflecting the spatial co-activation structure across electrodes. A vectorized representation of  $\mathcal{E}$  is denoted as  $\text{vec}(\mathcal{E})$ , a column vector of dimension  $\mathcal{V}^2$ .

The geodesic distance between two SPD matrices  $\mathcal{E}_1$  and  $\mathcal{E}_2$  is the same as the distance between their corresponding Cholesky matrices  $\mathcal{L}_1$  and  $\mathcal{L}_2$  (Lin, 2019) and is calculated as

$$d(\mathcal{L}_1, \mathcal{L}_2) = \left\{ \|\lfloor \mathcal{L}_1 \rfloor - \lfloor \mathcal{L}_2 \rfloor\|_F^2 + \|\log \mathbb{D}(\mathcal{L}_1) - \log \mathbb{D}(\mathcal{L}_2)\|_F^2 \right\}^{1/2}, \quad (1)$$

where  $\|\cdot\|_F$  denotes the Frobenius norm. Here,  $\mathcal{L}_1$  and  $\mathcal{L}_2$  are the *Cholesky factors* of the SPD matrices  $\mathcal{E}_1$  and  $\mathcal{E}_2$ , i.e., lower triangular matrices such that  $\mathcal{E} = \mathcal{L}\mathcal{L}^\top$ .

**EMG spectrograms:** for an EMG signal  $E_{\mathcal{V} \times \tau}$  collected from  $\mathcal{V}$  sensor nodes at a sampling frequency  $f_s$ , we compute the short-time Fourier transform (STFT) over successive time windows to obtain a power spectrogram representation  $\mathcal{S}_{\mathcal{V} \times F \times \tau'} = |\text{STFT}(E_{\mathcal{V} \times \tau})|^2$ , where  $F$  denotes the number of frequency bins and  $\tau'$  the number of time frames. Each slice  $\mathcal{S}_{\mathcal{V} \times F}^{(t)}$  captures the frequency-domain energy distribution of EMG activity across  $\mathcal{V}$  electrodes at time frame  $t$ . To reduce the spectral granularity, we bin the frequency axis into  $B$  frequency bands using  $\mathcal{B}_{\mathcal{V} \times B \times \tau'}(b) = \frac{1}{|F_b|} \sum_{f \in F_b} \mathcal{S}_{\mathcal{V} \times f \times \tau'}$ , where  $F_b$  is the set of frequency bins assigned to band  $b$ . In practice, we use either five bands  $B_1 = [80, 125]$  Hz,  $B_2 = [125, 250]$  Hz,  $B_3 = [250, 375]$  Hz,  $B_4 = [375, 687.5]$  Hz, and  $B_5 = [687.5, 1000]$  Hz, following Kaifosh et al. (2025), or 31 linearly spaced frequency bands between 80 and 1000 Hz. A vectorized representation of  $\mathcal{B}$  is denoted as  $\text{vec}(\mathcal{B})$ , a column vector of dimension  $\mathcal{V}B$ .

## 4.2 Audio (A)

**Audio spectrograms:** for a speech waveform  $a(t)$  sampled at frequency  $f_s$ , we compute a mel-scaled power spectrogram using a Hann-windowed short-time Fourier transform (STFT), followed by projection onto a mel filterbank with  $B$  mel bands. Specifically, we first obtain the power spectrogram  $\mathcal{M}_{F \times \tau'} = |\text{STFT}(a(t))|^2$ , where  $F$  denotes the number of frequency bins and  $\tau'$  the number of

time frames. This spectrogram is then projected onto a mel filterbank  $W_{\text{mel}}$  spanning the frequency range  $[f_{\min}, f_{\max}]$ , yielding

$$\mathcal{A}_{B \times \tau'}(b, t) = \sum_f W_{\text{mel}}(b, f) \mathcal{M}_{f, t}, \quad (331)$$

where  $b \in \{1, \dots, B\}$  indexes the mel bands. Each vector  $\mathcal{A}_B^{(t)}$  encodes the mel-band power distribution of the speech signal at frame  $t$ , emphasizing perceptually relevant frequency regions. We use  $B = 80$  mel bands,  $f_{\min} = 20$  Hz, and  $f_{\max} = f_s/2$ . We denote the column vector of an audio spectrogram by  $\mathcal{A}$  throughout the article.

**Audio features from SS models:** for a speech waveform  $a(t)$ , we extract self-supervised (SS) representations by passing the signal through a pre-trained model  $\mathcal{S}$ , yielding  $\mathcal{H} = \mathcal{S}(a(t))$ . The model  $\mathcal{S}$  can be instantiated as WAV2VEC 2.0 (Baevski et al., 2020), HUBERT (Hsu et al., 2021), or WAVLM (Chen et al., 2022). We denote the column vector of SS audio representations by  $\mathcal{H}$  throughout the article.

## 4.3 Sequence-to-sequence models

We construct sequences of  $\text{vec}(\mathcal{E})$ ,  $\text{vec}(\mathcal{B})$ ,  $\mathcal{A}$ , and  $\mathcal{H}$ , emitted every 20 ms and use a context length of 25 ms. For temporal relation modeling, we employ a causal time depth separable convolutional network (TDS), as described below.

We adapt the TDS model originally designed for EMG-based keyboard typing in Sivakumar et al. (2024) with minor modifications. The model relies exclusively on local temporal context, with a 1 s causal receptive field. To improve robustness to spatial variability in electrode activity, the architecture incorporates a *rotation-invariance* module consisting of a linear layer followed by a ReLU activation. This module is applied to electrode channel shifts of  $-1$ ,  $0$ , and  $+1$  positions, and the resulting outputs are averaged. The concatenated outputs from the rotation-invariance module are then fed into the TDS network for temporal modeling.

## 5 Results

### 5.1 $\mathcal{E}$ and $\mathbb{D}(\mathcal{E})$ encode articulatory information

Here we test whether covariance-based EMG features preserve discriminative structure related to articulation. We evaluate this on

DATA<sub>OROFACIAL GESTURES</sub>, where each trial is an orofacial movement recorded from 22 electrodes over 1.5 s. Each trial is represented by an EMG signal matrix  $E \in \mathbb{R}^{22 \times 7500}$ . We summarize each trial with a symmetric positive definite (SPD) covariance matrix  $\mathcal{E} \in \mathbb{R}^{22 \times 22}$ , and additionally consider its diagonal  $\mathbb{D}(\mathcal{E}) \in \mathbb{R}^{22}$ , which represents per-channel EMG power.

The vectors  $\mathbb{D}(\mathcal{E})$  corresponding to different orofacial gestures naturally form distinct clusters, as shown in figure 1. We further quantify this separability using the unsupervised  $k$ -medoids clustering algorithm (Kaufman and Rousseeuw, 1990), achieving an accuracy of 61.41% using  $\mathbb{D}(\mathcal{E})$  (averaged across 12 subjects). When using the full covariance matrix  $\mathcal{E}$  with the geodesic distance defined in equation 1, the  $k$ -medoids accuracy increases to 73.7%; both results are well above the random-chance level of 7.69%.

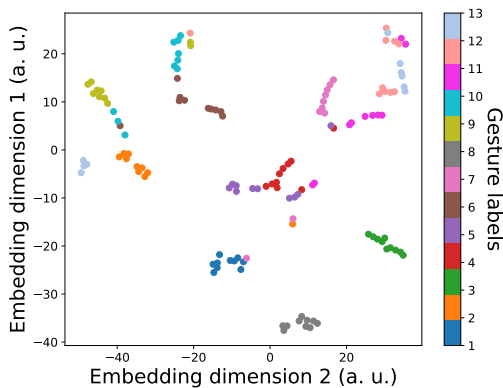


Figure 1: Different orofacial gestures are naturally separable using covariance-based EMG features.  $t$ -SNE visualization of diagonal vectors  $\mathbb{D}(\mathcal{E})$  for 13 orofacial movements from a single subject. Embedding is color-coded by articulatory gesture type ( $a.u.$  = arbitrary units).

These results demonstrate that both  $\mathcal{E}$  and  $\mathbb{D}(\mathcal{E})$  encode discriminative articulatory information. While  $\mathbb{D}(\mathcal{E})$  alone is sufficient to distinguish among different orofacial movements, incorporating the full covariance structure in  $\mathcal{E}$  improves decoding accuracy.

Note that other widely used EMG features, such as log-spectrograms (Sivakumar et al., 2024) or rectified time-domain signals (Halliday and Farmer, 2010), are not as straightforward to probe for this type of global structure. When raw EMG signals  $E \in \mathbb{R}^{22 \times 7500}$  are featurized using spectrograms or rectified signals, the temporal dimension may be reduced in granularity but is not collapsed into a sin-

gle frame. In contrast, covariance-based representations aggregate temporal information into a single fixed-dimensional feature, yielding  $\mathbb{D}(\mathcal{E}) \in \mathbb{R}^{22}$  or  $\mathcal{E} \in \mathbb{R}^{22 \times 22}$ . We analyze  $\mathbb{D}(\mathcal{E})$  using Euclidean distance, while  $\mathcal{E}$  is compared using the metric defined in equation 1. Consequently, commonly used time-frequency or time-domain features do not yield a directly comparable fixed-dimensional representation that captures articulatory structure in the same way.

## 5.2 $\mathcal{H}$ can linearly map to $\mathbb{D}(\mathcal{E})$

We test whether there exists a linear mapping defined by a weight matrix  $W$  and bias  $b$  such that  $\mathbb{D}(\mathcal{E}) \approx W\mathcal{H} + b$  with a high correlation<sup>2</sup>.

We use the training set described in section 3 DATA<sub>GENERAL</sub> to learn this mapping and evaluate it on the corresponding test set. We report the Pearson correlation between the predicted sequences  $\mathbb{D}(\mathcal{E}')$  and the ground-truth  $\mathbb{D}(\mathcal{E})$  on the test set. The representations  $\mathcal{H}$  are extracted using HUBERT (Hsu et al., 2021), WAV2VEC 2.0 (Baevski et al., 2020), and WAVLM (Chen et al., 2022). We evaluate BASE models with latent space dimension of 768 and 12 transformer layers, LARGE models with latent space dimension of 1024 and 24 transformer layers, and FINE-TUNED (FT) models that have been tuned for automatic speech recognition (ASR).

Correlation coefficients ( $r$ ) across models and layers are shown in figure 2. We find that a simple linear model can predict  $\mathbb{D}(\mathcal{E})$  from  $\mathcal{H}$  with a correlation as high as  $r = 0.85$ . The layer-wise trends across different models partially mirror the observations reported in (Cho et al., 2023, 2024) for electromagnetic articulography (EMA), where two local peaks were consistently observed across models. In our case, we observe two local peaks for WAV2VEC 2.0 models but only a single dominant peak for HUBERT and WAVLM models. A sharp decline in correlation emerges in the upper layers of fine-tuned models, reflecting the growing influence of task-specific objectives. This effect

<sup>2</sup>We actually aim to probe whether  $\mathcal{H}$  (768-1024 dimensions) can map to  $\text{vec}(\mathcal{E})$  (961 dimensions). However, the resulting  $\sim 10^6$ -parameter linear transformation would be severely ill-posed. To make this analysis tractable, we use  $\mathbb{D}(\mathcal{E})$  as a proxy because it provides a compact, well-conditioned, and physically meaningful representation grounded in articulatory mechanisms, making it well suited for linear probing. Importantly, this substitution is justified because both  $\mathcal{E}$  and  $\mathbb{D}(\mathcal{E})$  encode structured articulatory information, and the latter serves as a low-dimensional surrogate for the former, as shown in section 5.1.

is especially pronounced for WAV2VEC 2.0 compared to HUBERT and WAVLM.

Notably, for the HUBERT-BASE model, the peak correlation at layer 6 aligns with the layer previously identified as optimal for discrete speech resynthesis and spoken language modeling (Lakhotia et al., 2021). While prior work established this empirical result, the mechanistic basis for this peak remained unclear. Our analysis provides a principled interpretation: layer 6 exhibits the strongest linear predictive power for  $\mathbb{D}(\mathcal{E})$ , which encodes structured and discriminative articulatory information (i.e., different articulatory gestures such as tongue and jaw positions naturally form separable clusters). This tight alignment between articulatory structure and model representations offers a direct explanation for why layer 6 is particularly effective for downstream speech resynthesis and language modeling. In short, the layer that best captures articulatory mechanisms is also the one that yields the strongest downstream performance, providing convergent evidence for its functional role.

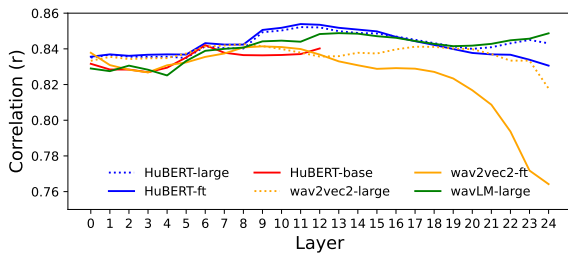


Figure 2: Layer-wise correlation ( $r$ ) between  $\mathbb{D}(\mathcal{E}')$  and  $\mathbb{D}(\mathcal{E})$  across different self-supervised speech models. A simple linear mapping is used to predict  $\mathbb{D}(\mathcal{E}')$  from  $\mathcal{H}$ .

We also examine whether a similar linear mapping exists between EMG spectrogram features ( $\text{vec}(\mathcal{B})$ ) and  $\mathcal{H}$ . Frequency bands of  $\mathcal{B}$  are obtained using five frequency bins, as described in section 4. However, the resulting correlation coefficients are substantially lower, with a maximum correlation of approximately  $r = 0.57$  (figure 3). For comparison, we also compute correlations for linear mapping between  $\mathcal{A}$  and  $\mathbb{D}(\mathcal{E})$  ( $r = 0.61$ ), which is considerably lower than the correlation between  $\mathcal{H}$  and  $\mathbb{D}(\mathcal{E})$ .

The above observations indicate that among the different EMG feature representations considered,  $\mathbb{D}(\mathcal{E})$  exhibits the strongest linear alignment with the self-supervised speech feature space  $\mathcal{H}$ . This strong correspondence suggests that  $\mathbb{D}(\mathcal{E})$  (consequently,  $\text{vec}(\mathcal{E})$ ) and  $\mathcal{H}$  encode highly compati-

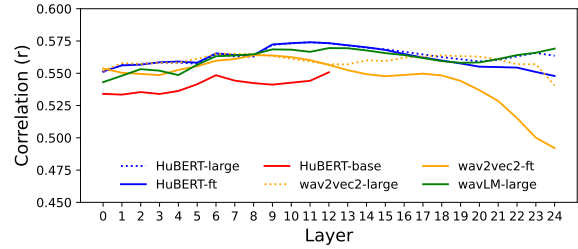


Figure 3: Layer-wise correlation ( $r$ ) between  $\mathcal{B}'$  and  $\mathcal{B}$  across different self-supervised speech models. A simple linear mapping is used to predict  $\mathcal{B}'$  from  $\mathcal{H}$ .

ble representations, making them particularly well suited for EMG-to-audio learning. In contrast, EMG spectrogram features ( $\mathcal{B}$ ) and their alignment with ( $\mathcal{H}$ ) yield notably weaker correlations.

### 5.3 emg2speech synthesis

As shown earlier, the following relationship holds:

$$\mathcal{H} \xrightarrow{\text{linear mapping}} \mathbb{D}(\mathcal{E}) \xrightarrow{\text{gesture-specific clustering}} \text{OROFACIAL MOVEMENTS.}$$

The existence of a simple linear mapping from  $\mathcal{H}$  to  $\mathbb{D}(\mathcal{E})$  is informative because it indicates that the self-supervised representations  $\mathcal{H}$  encode articulatory structure consistent with underlying muscle activations. This forward mapping is well posed:  $\mathcal{H}$  has moderate dimensionality (768–1024), whereas  $\mathbb{D}(\mathcal{E})$  is low dimensional (31), allowing the mapping to be estimated stably using linear regression.

By contrast, the inverse direction  $\mathbb{D}(\mathcal{E}) \rightarrow \mathcal{H}$  is intrinsically underdetermined and not uniquely invertible in the linear setting, since multiple high-dimensional speech representations can correspond to the same low-dimensional articulatory configuration. This challenge is further compounded when temporal alignment between EMG and audio is unknown. Nevertheless, the existence of a reliable forward mapping suggests that recovering  $\mathcal{H}$  from EMG is a feasible learning problem when using an appropriate nonlinear mapping<sup>3</sup>.

Motivated by this observation, we consider alignment-free prediction of  $\mathcal{H}$ -derived representations from EMG features ( $\text{vec}(\mathcal{E})$ ,  $\mathbb{D}(\mathcal{E})$ , or  $\text{vec}(\mathcal{B})$ ). Because the inverse linear mapping is ill posed, we model it using a nonlinear sequence-to-sequence architecture that can capture tempo-

<sup>3</sup>For linear probing, we use low-dimensional versions of both covariance and spectrogram representations:  $\mathbb{D}(\mathcal{E})$  and a 5-bin spectrogram  $\mathcal{B}$ . For speech synthesis, we instead use full-resolution features ( $\text{vec}(\mathcal{E})$  and a 31-bin  $\mathcal{B}$ ) to preserve fine-grained temporal and spectral structure.

ral and contextual dependencies present in  $\mathcal{H}$ . Concretely, EMG features are provided as input to a TDS convolutional network (section 4.3) that predicts discrete units associated with  $\mathcal{H}$ . We use the 100-unit discretization from layer 6 of HUBERT-BASE (Lakhotia et al., 2021), denoted  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$ . Training is performed with the connectionist temporal classification (CTC) loss (Graves et al., 2006), which enables learning without explicit frame-level alignment between EMG sequences and  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$ . Finally, the predicted  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  sequence is converted to audio using a pretrained Tacotron vocoder (Wang et al., 2017). The overall architecture is illustrated in figure 4.

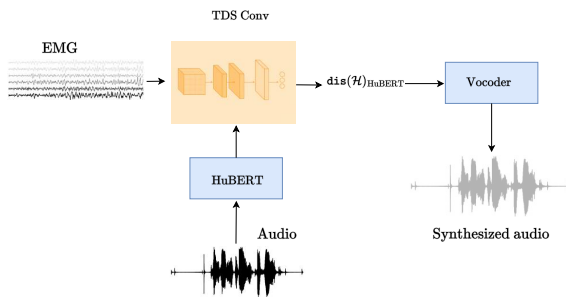


Figure 4: Multivariate EMG signals are converted into  $\text{vec}(\mathcal{E})$ ,  $\mathbb{D}(\mathcal{E})$ , or  $\mathcal{B}$ , and then passed through a TDS CONV block to predict  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$ , which are subsequently fed into a vocoder to synthesize audio. Frozen neural network components are shown in blue, and trainable components are shown in orange.

### 5.3.1 Results for $\text{DATA}_{\text{GENERAL}}$

We use Google text-to-speech (gTTS) to synthesize audio from the corresponding text transcripts. From this synthesized audio, we extract the discrete HuBERT units  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$ <sup>4</sup>. We report  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  decoding results in table 1.

We provide  $\text{vec}(\mathcal{E})$ ,  $\mathbb{D}(\mathcal{E})$ , or  $\text{vec}(\mathcal{B})$  as input to the TDS network, and train it to predict the corresponding  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  unit sequence using the CTC loss. For example, for the sentence  $\text{T-START} < \text{IT WAS PAID FOR} >_{\text{T-END}}$  with target

<sup>4</sup>For both  $\text{DATA}_{\text{GENERAL}}$  and  $\text{DATA}_{\text{ALS}}$ , we do not use subject-recorded audio for EMG-to-speech synthesis. The healthy participant vocalized the sentences during recording, whereas the ALS participant articulated them silently. In both cases, we rely only on transcript-based gTTS audio to derive  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$ . Subject-recorded audio from  $\text{DATA}_{\text{GENERAL}}$  is used only for probing linearity in the previous section. This design is motivated by clinical scenarios in which parallel EMG and audio recordings may be unavailable. Note that the gTTS audio is not temporally aligned with the EMG, which makes the translation problem more challenging than settings with paired, time-synchronized supervision.

units 71-12-71-12-4-12-4-40-93-86-13-58-32-1-99-..., the model learns a mapping from the EMG feature sequence to the target unit sequence. During inference, the model outputs a distribution over all 100  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  units at each time step, and we decode the most likely unit sequence using greedy search. For instance, the decoded sequence might be 71-12-57-4-54-40-93-86-13-58-16-14-76-6-36-.... We compute the unit error rate (UER) as the Levenshtein distance between the target and predicted  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  sequences, normalized by the target sequence length (see appendix D for ablation studies on training data size).

Table 1: Unit error rate (UER) for different EMG feature representations when predicting  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  units on  $\text{DATA}_{\text{GENERAL}}$ . The dataset and preprocessing details are described in section 3. Lower UER is better. Values are averaged over 5 random seeds.

MODEL INPUT	UER (% ↓)	INPUT FEATURE DIMENSION
$\text{vec}(\mathcal{B})$	$62.71 \pm 0.50$	961
$\mathbb{D}(\mathcal{E})$	$62.16 \pm 0.50$	31
$\text{vec}(\mathcal{E})$	<b><math>56.08 \pm 0.91</math></b>	961

We also present the results of phoneme-level decoding in table 2. For the sentence  $\text{T-START} < \text{IT WAS PAID FOR} >_{\text{T-END}}$  with the corresponding phonemic transcription  $\text{IH-T SPACE W-AA-Z SPACE P-EY-D SPACE F-AO-R}$ , the TDS model is trained to learn the mapping from  $\text{vec}(\mathcal{E})$ ,  $\mathbb{D}(\mathcal{E})$ , or  $\text{vec}(\mathcal{B})$  to phoneme sequences using the CTC loss. During inference, the model outputs probabilities for all 40 English phonemes at each time step, and the predictions are decoded using greedy search. For example, the decoded output might be  $\text{IH-T SPACE W-AA-Z SPACE P-EY-T SPACE F-AO-R}$ . We compute the phoneme error rate (PER) as the Levenshtein distance between the target and decoded phoneme sequences, normalized by the length of the target sequence.

Table 2: Phoneme error rate (PER) for different EMG feature representations when predicting phonemes on  $\text{DATA}_{\text{GENERAL}}$ . The dataset and preprocessing details are described in section 3. Lower PER is better. Values are averaged over 5 random seeds.

MODEL INPUT	PER (% ↓)	INPUT FEATURE DIMENSION
$\text{vec}(\mathcal{B})$	$44.40 \pm 2.28$	961
$\mathbb{D}(\mathcal{E})$	$44.40 \pm 1.51$	31
$\text{vec}(\mathcal{E})$	<b><math>32.78 \pm 0.66</math></b>	961

As shown in tables 1 and 2,  $\text{vec}(\mathcal{E})$  outperforms  $\text{vec}(\mathcal{B})$ .  $\mathcal{B}$  was computed using 31 linearly spaced frequency bins, and for any given time frame, both  $\text{vec}(\mathcal{E})$  and  $\text{vec}(\mathcal{B})$  have 961 dimensions. Notably, even  $\mathbb{D}(\mathcal{E})$ , which has only 31 dimensions (i.e.,  $31 \times$  fewer dimensions than  $\text{vec}(\mathcal{B})$ ), performs on par with  $\text{vec}(\mathcal{B})$ . This finding is consistent with the linear mapping results shown in figures 2 and 3.

**PHONEME GUIDED DECODING OF  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$ :** As shown in table 2, phoneme sequences can be decoded more accurately than  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  unit sequences. Motivated by this observation, we train the TDS convolutional encoder in figure 4 with two prediction heads: one for phonemes, producing framewise posteriors  $P(\text{PHONEME} \mid \text{EMG})$ , and one for  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  units, producing framewise posteriors  $P(\text{dis}(\mathcal{H})_{\text{HUBERT}} \mid \text{EMG})$ . The model is optimized with CTC losses for both outputs, together with an additional consistency loss that encourages phoneme-consistent  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  predictions. Specifically, we use a precomputed lookup table  $P(\text{PHONEME} \mid u)$ , where  $u \in \mathcal{U}$  and  $\mathcal{U} = \text{dis}(\mathcal{H})_{\text{HUBERT}}$  denotes the discrete HuBERT unit set, to transform the unit posterior at each frame  $t$  into a phoneme distribution by marginalizing over units:

$$\tilde{P}(\text{PHONEME} \mid \text{EMG})_t = \sum_{u \in \mathcal{U}} P(\text{PHONEME} \mid u) P(u \mid \text{EMG})_t.$$

We then minimize a cross-entropy loss between  $\tilde{P}(\text{PHONEME} \mid \text{EMG})_t$  and the phoneme-head posterior  $P(\text{PHONEME} \mid \text{EMG})_t$  (see appendix B for more details). At inference time, we decode only  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  units from the unit head using greedy decoding. Resulting improvement is shown in table 3.

Table 3: Effect of phoneme-guided training on unit decoding with  $\text{vec}(\mathcal{E})$  as input on  $\text{DATA}_{\text{GENERAL}}$ . Values are averaged over 5 random seeds. The improvement is statistically significant ( $p < 10^{-6}$ ).

TRAINING OBJECTIVE	UER (% ↓)
Unit CTC only	56.08 ± 0.91
Phoneme-guided decoding	<b>51.81 ± 0.62</b>

Furthermore, three human raters (see appendix D.3) listened to all 400 synthesized audio samples in the test set (generated from  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  units obtained via phoneme-guided decoding with

$\text{vec}(\mathcal{E})$  as input) and transcribed them. We compute the word error rate (WER) as the Levenshtein distance between each rater’s transcription and the ground-truth transcript, normalized by the length of the ground-truth transcript. We report the resulting WERs in table 5.

### 5.3.2 Results for $\text{DATA}_{\text{ALS}}$



We follow the same preprocessing, model architecture, and training procedure used for  $\text{DATA}_{\text{GENERAL}}$ . We report unit decoding performance in table 4. To assess end-to-end intelligibility, we synthesize speech from the decoded units and measure word error rate (WER) using human transcriptions (on all 60 synthesized audios in the test set); the resulting WER is reported in table 5.

Table 4: Unit error rate (UER) for different EMG feature representations when predicting  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  units on  $\text{DATA}_{\text{ALS}}$ . Dataset and preprocessing details are described in section 3. Lower UER is better. Values are averaged over 5 random seeds.

MODEL INPUT	UER (% ↓)
$\text{vec}(\mathcal{B})$	59.18 ± 2.81
$\mathbb{D}(\mathcal{E})$	54.27 ± 0.56
$\text{vec}(\mathcal{E})$	52.29 ± 0.91
$\text{vec}(\mathcal{E})$ with phoneme-guided decoding	<b>46.98 ± 1.11</b>

Table 5: WER as rated by human transcribers.

HUMAN TRANSCRIBER	$\text{DATA}_{\text{GENERAL}}$ WER (% ↓)	$\text{DATA}_{\text{ALS}}$ WER (% ↓)
1	65.09	52.69
2	59.32	50.97
3	63.23	48.81
Average	<b>62.55 ± 2.40</b>	<b>50.82 ± 1.59</b>

We contextualize these WERs relative to prior brain-computer interface studies in appendix C. Please see the following links for demonstrations:  [Audio](#). We also provide the ground-truth transcripts and test-set transcripts from three human raters:  [Transcriptions](#).

## 6 Conclusion

We present methods and datasets that convert orofacial EMG signals directly into speech, and we demonstrate the system with both a healthy participant and a participant diagnosed with ALS.

## 7 Limitations

This work primarily focuses on the technical aspects of EMG-to-speech modeling, including characterizing the structure of orofacial EMG signals, quantifying their relationship to self-supervised speech representations, and designing encoder architectures grounded in articulatory mechanisms. Our clinical demonstration uses approximately one hour of data from a single participant with ALS. As a result, we do not yet characterize how performance evolves under day-to-day distribution shifts in EMG signals (e.g., changes in electrode placement, skin impedance, fatigue, or disease progression). Consequently, we also do not evaluate whether this non-stationarity is more or less challenging than the distribution shifts observed in invasive neural interfaces (Wairagkar et al., 2025; Willett et al., 2023).

In addition, we do not demonstrate sustained, long-term performance of this non-invasive neuroprosthesis. In contrast, prior work on invasive neuroprostheses has reported stability over extended periods in related decoding settings, including brain-to-text (Fan et al., 2023) and cursor-based brain-computer interfaces (Wilson et al., 2025).

Finally, we do not explore whether large-scale pretrained EMG-to-speech models can improve decoding performance. For example, in related work on EMG-based keyboard typing (EMG2QWERTY) (Sivakumar et al., 2024), pretraining on data from 100 individuals improved accuracy after fine-tuning to new individuals, although zero-shot performance remained limited. Speech may be more challenging than discrete key typing, and future work should investigate how to build and effectively leverage large-scale pretrained models for EMG-to-speech translation.

We are actively addressing these limitations through ongoing longitudinal studies and by expanding data collection to build larger EMG-to-speech corpora from individuals with diverse clinical etiologies, including ALS and laryngectomy.

## 8 Ethical considerations

Research was conducted in accordance with the principles embodied in the Declaration of Helsinki and with approval from the home institution’s Institutional Review Board (IRB). All participants provided written informed consent. All participants also provided consent for publication of deidentified data. Volunteers of any gender and from

all racial and ethnic groups were eligible to participate. Participants were required to be at least 18 years old, able to understand spoken and written English, and able to follow task instructions. Participants had no skin conditions or wounds at electrode placement sites and were excluded if they had uncorrected vision problems. Children, individuals unable to provide informed consent, and prisoners were not included in the experiments.

The participant with ALS was first diagnosed in 2019 and has non-familial ALS with spasticity.

## References

- Alexei Baeviski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. 2020. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460.
- Tyler Benster, Guy Wilson, Reshef Elisha, Francis R Willett, and Shaul Druckmann. 2024. A cross-modal approach to silent speech with llm-enhanced recognition. *arXiv preprint arXiv:2403.05583*.
- Sanyuan Chen, Chengyi Wang, Zhengyang Chen, Yu Wu, Shujie Liu, Zhuo Chen, Jinyu Li, Naoyuki Kanda, Takuya Yoshioka, Xiong Xiao, and 1 others. 2022. Wavlm: Large-scale self-supervised pre-training for full stack speech processing. *IEEE Journal of Selected Topics in Signal Processing*, 16(6):1505–1518.
- Cheol Jun Cho, Abdelrahman Mohamed, Alan W Black, and Gopala K Anumanchipalli. 2024. Self-supervised models of speech infer universal articulatory kinematics. In *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 12061–12065. IEEE.
- Cheol Jun Cho, Peter Wu, Abdelrahman Mohamed, and Gopala K Anumanchipalli. 2023. Evidence of vocal tract articulation in self-supervised learning of speech. In *ICASSP 2023-2023 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 1–5. IEEE.
- Lorenz Diener, Gerrit Felsch, Miguel Angrick, and Tanja Schultz. 2018. Session-independent array-based emg-to-speech conversion using convolutional neural networks. In *Speech Communication; 13th ITG-Symposium*, pages 1–5.
- Chaofei Fan, Nick Hahn, Foram Kamdar, Donald Avansino, Guy Wilson, Leigh Hochberg, Krishna V Shenoy, Jaimie Henderson, and Francis Willett. 2023. Plug-and-play stability for intracortical brain-computer interfaces: a one-year demonstration of seamless brain-to-text communication. *Advances in neural information processing systems*, 36:42258–42270.



855	Viswanath Sivakumar, Jeffrey Seely, Alan Du, Sean R Bittner, Adam Berenzweig, Anuoluwapo Bolarinwa, Alexandre Gramfort, and Michael I Mandel. 2024. emg2qwerty: A large dataset with baselines for touch typing using surface electromyography. In <i>The Thirty-eighth Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	908
856		909
857		910
858		911
859		912
860		913
861		914
862	Arthur R. Toth, Michael Wand, and Tanja Schultz. 2009. Synthesizing speech from electromyography using voice transformation techniques. In <i>Interspeech 2009</i> , pages 652–655.	915
863		916
864		917
865		918
866	Maitreyee Wairagkar, Nicholas S Card, Tyler Singer-Clark, Xianda Hou, Carrina Iacobacci, Lee M Miller, Leigh R Hochberg, David M Brandman, and Sergey D Stavisky. 2025. An instantaneous voice-synthesis neuroprosthesis. <i>Nature</i> , pages 1–8.	919
867		920
868		921
869		922
870		923
871	Yuxuan Wang, RJ Skerry-Ryan, Daisy Stanton, Yonghui Wu, Ron J. Weiss, Navdeep Jaitly, Zongheng Yang, Ying Xiao, Zhifeng Chen, Samy Bengio, Quoc Le, Yannis Agiomyrgiannakis, Rob Clark, and Rif A. Saurous. 2017. Tacotron: Towards end-to-end speech synthesis.	924
872		925
873		926
874		927
875		928
876		929
877	Francis R Willett, Erin M Kunz, Chaofei Fan, Donald T Avansino, Guy H Wilson, Eun Young Choi, Foram Kamdar, Matthew F Glasser, Leigh R Hochberg, Shaul Druckmann, and 1 others. 2023. A high-performance speech neuroprosthesis. <i>Nature</i> , 620(7976):1031–1036.	930
878		931
879		932
880		933
881		934
882		935
883	Guy H Wilson, Elias A Stein, Foram Kamdar, Donald T Avansino, Tsam Kiu Pun, Ronnie Gross, Tommy Hosman, Tyler Singer-Clark, Anastasia Kapitonava, Leigh R Hochberg, and 1 others. 2025. Long-term unsupervised recalibration of cursor-based intracortical brain–computer interfaces using a hidden markov model. <i>Nature Biomedical Engineering</i> , pages 1–19.	936
884		937
885		938
886		939
887		940
888		941
889		942
890	Maria K Wolters, Karl B Isaac, and Steve Renals. 2010. Evaluating speech synthesis intelligibility using amazon mechanical turk. In <i>Proc. 7th Speech Synthesis Workshop (SSW7)</i> , pages 136–141.	943
891		944
892		945
893		946

## 894 A Experimental details

895 We collect EMG signals from 31 sites on the neck,  
896 chin, jaw, cheek, and lips using monopolar elec-  
897 trodes. An ACTICHAMP PLUS amplifier and asso-  
898 ciated active electrodes from BRAINVISION ([Brain](#)  
899 [Vision](#)) are used to record EMG signals at 5000 Hz.  
900 To ensure proper contact between the skin sur-  
901 face and electrodes, we use SUPERVISC, a high-  
902 viscosity electrolyte gel from EASYCAP ([Easycap](#)).  
903 We develop a software suite in a PYTHON environ-  
904 ment to provide visual cues to participants and to  
905 collate and store timestamped data. For time syn-  
906 chronization, we use Lab Streaming Layer ([LSL](#)).  
907 See figure 5 for electrode placement. In addition to

the 31 data electrodes, we also use a GROUND elec-  
trode (marked as GND) and a REFERENCE electrode  
(marked as 32). The GROUND electrode is placed  
on the left earlobe and the REFERENCE electrode  
is placed on the right earlobe.

Before signal acquisition, participants are  
briefed on the experimental protocol and seated  
comfortably in a chair. Sentence start and end  
times are timestamped using mouse clicks. When  
a participant is ready to articulate a sentence, they  
click the mouse to prompt the sentence to appear  
on the screen. Once articulation is complete, they  
click again to indicate the end, which causes the  
sentence to disappear. This allows participants to  
articulate at their own pace.

The data collection environment is carefully con-  
trolled to reduce AC electrical interference. EMG  
signals undergo minimal preprocessing. The sig-  
nal from the REFERENCE channel (electrode 32)  
is subtracted from all other channels. The result-  
ing signals are bandpass filtered using a third-order  
Butterworth filter between 80 and 1000 Hz and seg-  
mented according to sentence start and end times  
based on synchronized timestamps. The segmented  
sentences are subsequently  $z$ -normalized along the  
time dimension for each channel.

The electrodes are positioned over regions that  
overlay muscle groups involved in speech articula-  
tion, providing coverage of key articulators such  
as the tongue, jaw, lips, and larynx. Electrode lo-  
cations 19, 21, 3, and 1 approximately overlie the  
*hyoglossus*, *palatoglossus*, and *styloglossus* mus-  
cles. These muscles are located in the lower cheek  
region and play a vital role in tongue movement.  
They are also consistently recruited across a wide  
range of articulatory gestures. Muscles in the upper  
and posterior cheek regions include the *masseter*  
and *temporalis*, which control jaw motion, and the  
*zygomaticus*, which is involved in upper lip ele-  
vation. These muscles correspond approximately  
to electrode regions around nodes 22, 18, 17, and  
15 in figure 5. Electrodes located beneath the jaw  
capture activity from muscles involved in tongue  
protrusion and jaw-tongue coordination, such as  
the *genioglossus* near electrodes 8, 9, 23, and 25,  
as well as the *digastric*. Additionally, electrodes  
near the laryngeal region, including nodes 6, 7, 10,  
11, 26, and 27, reflect activity from muscles that  
modulate laryngeal and hyoid position, such as the  
*sternohyoid*, *stylohyoid*, and *digastric*. These mus-  
cles contribute to pitch control, vowel shaping, and  
jaw movement.

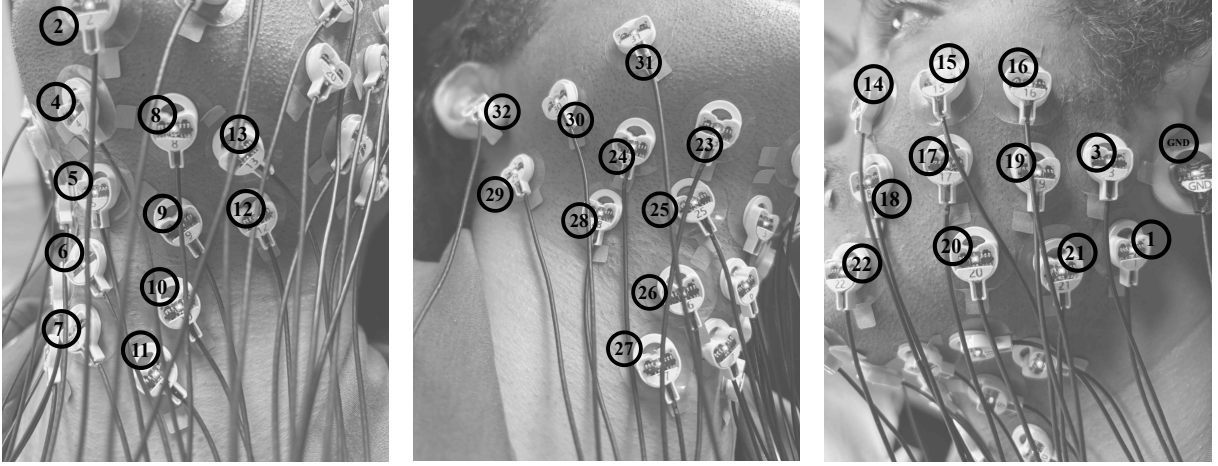


Figure 5: LEFT: Electrode placement on the left side of the neck. MIDDLE: Electrode placement on the right side of the neck. RIGHT: Electrode placement on the left cheek.

## B Detailed explanation: phoneme guided decoding of $\text{dis}(\mathcal{H})_{\text{HUBERT}}$

We train a bidirectional gated recurrent unit (GRU) model to map  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  unit sequences to phoneme sequences using a CTC objective on  $\text{DATA}_{\text{GENERAL}}$  (train-val-test split as described in section 3). After CTC decoding, the model achieves  $\text{PER} = 0\%$  (phoneme error rate) on the corresponding test split. Using this trained model, we estimate a unit-to-phoneme conditional table by aggregating framewise posteriors: for each HUBERT unit  $u$ , we collect the predicted phoneme distribution  $P_{\text{GRU}}(\text{PHONEME}_t | u)$  at every frame  $t$  where  $u_t = u$ , and average across all such frames, i.e.,

$$P(\text{PHONEME} | u) = \frac{1}{N_u} \sum_{t: u_t=u} P_{\text{GRU}}(\text{PHONEME}_t = \text{PHONEME} | u),$$

where  $N_u = |\{t : u_t = u\}|$  denotes the total number of frames in the dataset for which the HUBERT unit equals  $u$ . We remove the CTC blank symbol from  $P(\cdot | u)$  and renormalize. This yields  $P(\text{PHONEME} | \text{dis}(\mathcal{H})_{\text{HUBERT}})$ , which we use as a fixed probabilistic mapping in our consistency regularization.

For phoneme-guided decoding of  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$ , we train the TDS convolutional encoder in figure 4 with two heads: one for phonemes, producing framewise posteriors  $P(\text{PHONEME} | \text{EMG})$ , and one for  $\text{dis}(\mathcal{H})_{\text{HUBERT}}$  units, producing framewise posteriors  $P(u | \text{EMG})$ , where  $u \in \mathcal{U}$  and

$\mathcal{U} = \text{dis}(\mathcal{H})_{\text{HUBERT}}$  denotes the discrete HuBERT unit set. The model is optimized with CTC losses for both outputs. Additionally, we impose a consistency loss by using the precomputed lookup table  $P(\text{PHONEME} | u)$  to transform the unit posterior at each frame  $t$  into a phoneme distribution via marginalization over units:

$$\tilde{P}(\text{PHONEME} | \text{EMG})_t = \sum_{u \in \mathcal{U}} P(\text{PHONEME} | u) \cdot P(u | \text{EMG})_t.$$

We then minimize a cross-entropy loss between  $\tilde{P}(\text{PHONEME} | \text{EMG})_t$  and the phoneme-head posterior  $P(\text{PHONEME} | \text{EMG})_t$ :

$$\mathcal{L}_{\text{cons}} = - \sum_t \sum_{\text{PHONEME}} P(\text{PHONEME} | \text{EMG})_t \cdot \log \tilde{P}(\text{PHONEME} | \text{EMG})_t.$$

The total training objective is a weighted sum of the two CTC losses and the proposed consistency term:

$$\mathcal{L}_{\text{total}} = \lambda_{\text{unit}} \mathcal{L}_{\text{CTC}}^{\text{unit}} + \lambda_{\text{phone}} \mathcal{L}_{\text{CTC}}^{\text{phone}} + \lambda_{\text{cons}} \mathcal{L}_{\text{cons}}.$$

In our experiments, we use  $\lambda_{\text{unit}} = 0.8$ ,  $\lambda_{\text{phone}} = 0.1$ , and  $\lambda_{\text{cons}} = 0.1$ .

We further probe the structure of  $P(\text{PHONEME} | \text{dis}(\mathcal{H})_{\text{HUBERT}})$ . In figure 6, we visualize, for each HUBERT unit, the phoneme with the highest conditional probability (i.e.,  $\text{argmax}_{\text{PHONEME}} P(\text{PHONEME} | u)$ ). Multiple HUBERT units may map to the same most-probable phoneme  $p$ . In such cases, for each phoneme  $p$ , we order the corresponding units by increasing entropy of  $P(\cdot | u)$ ; lower entropy indicates a sharper,

1018 more confident association between the unit and  
1019 phoneme  $p$ . We find that alveolar consonants (e.g.,  
1020 T, S, N) and A and I based vowels (e.g., AA, AY,  
1021 AH, IY, IH) have many HUBERT units mapping to  
1022 them.

1023 This is consistent with prior analyses showing  
1024 that these phones are among the most frequently  
1025 occurring in conversational English and that their  
1026 phonetic realizations have different manifesta-  
1027 tions depending on coarticulatory context (Mines  
1028 et al., 1978). These observations suggest that our  
1029 lookup dictionary  $P(\text{PHONEME} \mid \text{dis}(\mathcal{H})_{\text{HUBERT}})$   
1030 captures meaningful phonetic structure and is  
1031 grounded in known articulatory regularities of  
1032 speech.

### 1033 C Detailed literature review

1034 Here, we review prior work on speech neural and  
1035 neuromuscular interfaces and contextualize our re-  
1036 sults relative to state-of-the-art methods. A substan-  
1037 tial body of research (Jou et al., 2006; Kapur et al.,  
1038 2020; Meltzner et al., 2018; Toth et al., 2009; Janke  
1039 and Diener, 2017; Diener et al., 2018; Littlejohn  
1040 et al., 2025) has laid the groundwork for EMG-  
1041 based speech interfaces. Among the earliest stud-  
1042 ies, Jou et al. (2006) demonstrate EMG-to-speech  
1043 conversion on a small corpus of 50 sentences. Ka-  
1044 pur et al. (2020) use a corpus of 15 sentences and,  
1045 rather than performing phoneme-level decoding,  
1046 formulate the task as a 15-way classification prob-  
1047 lem. Meltzner et al. (2018) study EMG-to-text  
1048 recognition for isolated words, phrases drawn from  
1049 a  $\sim 200$ -word vocabulary, and continuous sentences  
1050 using a custom grammar-based recognition model  
1051 over a set of 1200 scripted phrases. Toth et al.  
1052 (2009) present EMG-to-speech conversion on a  
1053 corpus of 500 sentences. Janke and Diener (2017)  
1054 demonstrate EMG-to-speech conversion using up  
1055 to two hours of data and 2000 utterances.

1056 Overall, these studies rely on private datasets  
1057 and task-specific pipelines, and they typically eval-  
1058 uate on small, constrained corpora. In addition,  
1059 the works do not release full implementations (e.g.,  
1060 code repositories) or sufficient methodological de-  
1061 tails to enable direct reproducibility. As a result, it  
1062 is difficult to directly compare performance across  
1063 systems, and all the above results do not establish  
1064 generalization to open-vocabulary English settings.

1065 A reproducible benchmark for open-vocabulary  
1066 EMG-to-speech conversion was introduced by  
1067 Gaddy and Klein (2020, 2021). However, these

1068 works rely on time-aligned EMG-audio pairs for  
1069 training. Building on Gaddy and Klein (2021), Ben-  
1070 ster et al. (2024) propose an approach that lever-  
1071 ages an audio-only corpus in addition to paired  
1072 EMG-audio data. While effective in the bench-  
1073 mark setting, such methods can be difficult to  
1074 deploy in clinical scenarios where parallel EMG-  
1075 audio recordings may be unavailable or unreliable.  
1076 On the large-vocabulary corpus, Gaddy and Klein  
1077 (2020) report a word error rate (WER) of 68%,  
1078 and Gaddy and Klein (2021) reduce this to 42%.  
1079 Benster et al. (2024) report WER  $< 10\%$  on the  
1080 Gaddy and Klein (2020) dataset by using a large  
1081 language model (LLM) to post-correct the inter-  
1082 mediate EMG-to-phoneme output. However, their  
1083 system no longer supports streaming synthesis, and  
1084 the evaluation does not fully characterize potential  
1085 information leakage through the LLM (e.g., mem-  
1086 orization or exposure to overlapping text distribu-  
1087 tions). Consequently, their results are not directly  
1088 comparable to strictly streaming EMG-to-speech  
1089 systems. Littlejohn et al. (2025) report a WER of  
1090 74% on the Gaddy and Klein (2020) dataset us-  
1091 ing a CNN+RNN transducer model; however, their  
1092 train-test splits and implementation details are not  
1093 publicly available, which prevents direct compar-  
1094 ison. In our setting, we address a harder learning  
1095 problem by not assuming time-aligned EMG-audio  
1096 pairs during training, and we report a WER of 62%  
1097 on an open-vocabulary corpus. We emphasize that  
1098 these WER values should not be compared one-to-  
1099 one across studies, since the data collection setup,  
1100 training targets and alignment assumptions, prob-  
1101 lem formulation, and evaluation methodology dif-  
1102 fer substantially. We report these results to provide  
1103 context relative to prior EMG-based speech inter-  
1104 faces.

1105 To address these limitations, we build on widely  
1106 available pretrained speech models and vocoders,  
1107 but adapt them to the EMG setting through princi-  
1108 pled, articulatorily motivated design choices. We  
1109 characterize the structure of orofacial EMG signals  
1110 and introduce methods that are straightforward to  
1111 implement yet carefully tailored for this problem,  
1112 yielding substantial gains. For example, we repre-  
1113 sent EMG signals using covariance matrices and  
1114 propose phoneme-guided decoding of speech units  
1115 (section 5.3.1), which exploits phonetic structure  
1116 to improve the fidelity of synthesized speech. Our  
1117 encoder design and phoneme-guided decoding are  
1118 grounded in articulatory mechanisms and estab-  
1119 lished regularities of speech production.

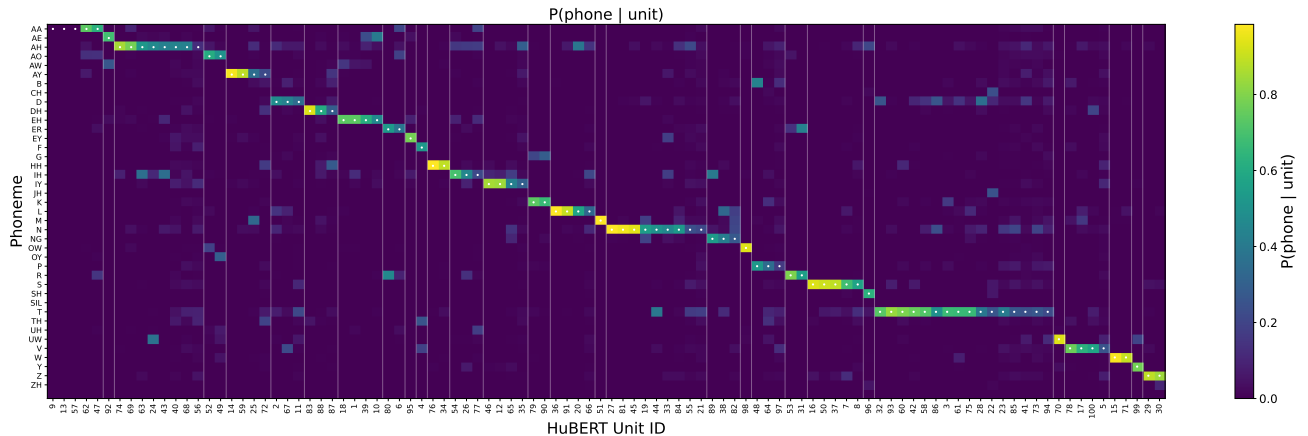


Figure 6: For each HUBERT unit, the phoneme with the highest conditional probability (i.e.,  $\operatorname{argmax}_{\text{PHONEME}} P(\text{PHONEME} | u)$ ) is shown. Multiple HUBERT units may map to the same most-probable phoneme  $p$ . In such cases, for each phoneme  $p$ , we order the corresponding units from left to right by increasing entropy of  $P(\cdot | u)$ .

1120 The WERs reported in this work fall broadly  
 1121 within the range observed for invasive speech neural  
 1122 interfaces. For example, [Wairagkar et al. \(2025\)](#)  
 1123 report a median WER of 43% using 256 intracortical  
 1124 electrodes, trained on approximately 8300 cued  
 1125 sentences spanning about 38 hours of data; the  
 1126 large number of hours for a comparable number of  
 1127 sentences reflects the substantially slower speak-  
 1128 ing rate in that study relative to ours. Similarly,  
 1129 [Metzger et al. \(2023\)](#) report a WER of 54% using  
 1130 253 electrodes with a 1024-word vocabulary. In  
 1131 our setting, we achieve a WER of 62% on a large-  
 1132 vocabulary corpus. These error rates should be in-  
 1133 terpreted in the context of neural speech interfaces,  
 1134 where the sensing modality, signal-to-noise ratio,  
 1135 data size, and evaluation setup differ substantially  
 1136 from conventional ASR (automatic speech  
 1137 recognition) benchmarks.

## 1138 D Additional technical details

### 1139 D.1 Effect of training data size

1140 We study how training-set size affects the unit error  
 1141 rate (UER) when decoding  $\operatorname{dis}(\mathcal{H})_{\text{HUBERT}}$ . Using  
 1142  $\text{DATA}_{\text{GENERAL}}$ , we train the model with between  
 1143 2000 and 8000 sentences, while keeping the valida-  
 1144 tion and test splits fixed as in section 3. Over this  
 1145 range, we observe an approximately linear improve-  
 1146 ment in UER with increasing training data. This  
 1147 differs from the power-law behavior commonly  
 1148 reported for large-scale from-scratch training ([Ka-  
 1149 plan et al., 2020](#)). Although we cannot draw strong  
 1150 conclusions from a single-participant study, these  
 1151 results suggest that, when leveraging frozen pre-

1152 trained speech representations, performance may  
 1153 remain data-limited and continue to benefit pre-  
 1154 dictably from additional EMG training data, yield-  
 1155 ing lower UER (and consequently WER) as the  
 1156 dataset scales.

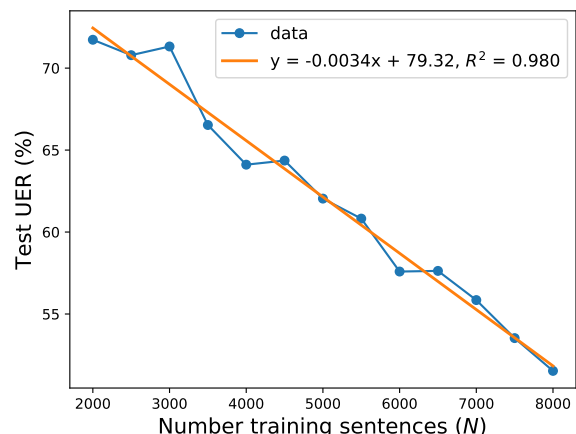


Figure 7: Test UER as a function of the number of training sentences. Over the evaluated range (2000-8000 sentences), UER decreases approximately linearly with increasing training data, indicating consistent gains from additional EMG data. UER is computed by decoding  $\operatorname{dis}(\mathcal{H})_{\text{HUBERT}}$  from a model trained with  $\operatorname{vec}(\mathcal{E})$  as input; during training, we use phoneme-guided decoding.

### 1157 D.2 Losses

1158 As shown in figure 8, all loss terms decrease  
 1159 smoothly during the early stages of training, indi-  
 1160 cating stable optimization. Validation losses begin  
 1161 to increase after approximately 25 epochs, suggest-  
 1162 ing the onset of overfitting. The consistency loss

1163  $\mathcal{L}_{\text{cons}}$  and the phoneme-level CTC loss  $\mathcal{L}_{\text{CTC}}^{\text{phone}}$  de- 1196  
 1164 crease more rapidly than the unit-level loss  $\mathcal{L}_{\text{CTC}}^{\text{unit}}$ , 1197  
 1165 consistent with their role as auxiliary objectives  
 1166 that regularize training and encourage better align-  
 1167 ment for unit prediction. Together with figure 7,  
 1168 which shows that UER decreases approximately  
 1169 linearly with training set size over the evaluated  
 1170 range and has not yet saturated, and the overfitting  
 1171 observed in figure 8, these results suggest that the  
 1172 model remains data-limited and can continue to  
 1173 benefit from additional training data.

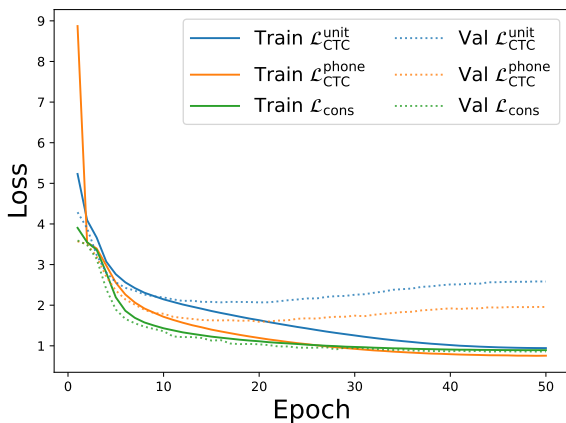


Figure 8: Training and validation losses versus epoch. The causal TDS convolutional model in section 4.3 is trained for 50 epochs with  $\text{vec}(\mathcal{E})$  as input, using  $\mathcal{L}_{\text{CTC}}^{\text{unit}}$ ,  $\mathcal{L}_{\text{CTC}}^{\text{phone}}$ , and  $\mathcal{L}_{\text{cons}}$ .

### 1174 D.3 Transcriptions

1175 Human evaluation of synthesized speech is com- 1196  
 1176 monly used to assess the intelligibility and overall 1197  
 1177 quality of generated audio (Wolters et al., 2010).  
 1178 For both  $\text{DATA}_{\text{GENERAL}}$  and  $\text{DATA}_{\text{ALS}}$ , we ask hu-  
 1179 man raters to listen to the synthesized audio and  
 1180 transcribe each utterance in English. Raters are  
 1181 not restricted to a predefined vocabulary and may  
 1182 write any English words for both corpora (even  
 1183 though  $\text{DATA}_{\text{ALS}}$  contains only about 300 unique  
 1184 words, raters are not informed of this constraint).  
 1185 In this sense, our evaluation targets open-vocabulary  
 1186 recognition and is less constrained than evalua-  
 1187 tions such as Metzger et al. (2023) and Littlejohn  
 1188 et al. (2025), which use fixed vocabularies of 1,024  
 1189 words.

1190 In figure 9, we summarize the distribution of  
 1191 WER and PER across three human transcribers  
 1192 for all evaluated utterances (460 sentences in total  
 1193 across  $\text{DATA}_{\text{GENERAL}}$  and  $\text{DATA}_{\text{ALS}}$ ). The raters  
 1194 exhibit similar central tendency and spread, indi-  
 1195 cating strong agreement. To compute PER, we

phonemize each rater’s transcription and compare  
 it against the ground-truth phonemized sentence.

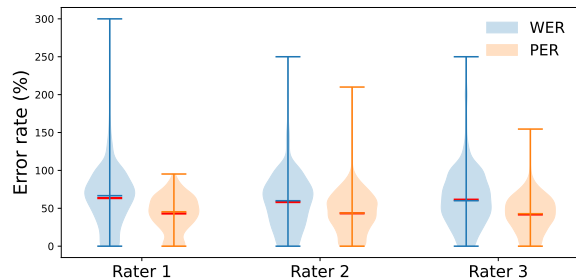


Figure 9: Distributions of WER and PER across three human transcribers for 460 utterances. Means are shown in red.

1198 Across raters, the mean PER is 42.79%, which  
 1199 is lower than the mean WER of 61.02%. This  
 1200 suggests that, even when transcribed words are in-  
 1201 correct, the errors are often phonetically plausible.

1202 Furthermore, we train the TDS convolutional  
 1203 model to decode phonemes on all 460 sentences  
 1204 (following the procedure in section 4.3) and then  
 1205 map the predicted phoneme sequences to words  
 1206 using a weighted finite-state transducer (WFST)  
 1207 decoder<sup>5</sup>.

1208 For phoneme-to-word decoding, we use the  
 1209 LibriSpeech-100 transcripts (Panayotov et al.,  
 1210 2015), which contain roughly 38000 sentences and  
 1211 35000 unique words. From these transcripts, we  
 1212 build a pronunciation lexicon WFST,  $L$ , that maps  
 1213 phoneme sequences to words. We also train a  
 1214 4-gram language model with KenLM (Heafield,  
 1215 2011) and convert it into a grammar WFST,  $G$ . Fi-  
 1216 nally, we construct the CTC topology WFST,  $H$ ,  
 1217 which encodes the allowable label sequences under  
 1218 the CTC criterion.

1219 We compose these components into a decoding  
 1220 graph (Mohri et al., 2008),  $HLG = H \circ L \circ G$ ,  
 1221 which integrates the CTC constraints ( $H$ ), lexicon  
 1222 mapping ( $L$ ), and language model probabilities ( $G$ ).  
 1223 At inference time, we perform beam search over  
 1224  $HLG$  with beam width 50 and compute WER as  
 1225 the normalized Levenshtein distance between the  
 1226 reference and decoded word sequences. This pro-  
 1227 cedure yields a WER of 51.17% and a PER of  
 1228 38.19%. The language model is trained only on  
 1229 LibriSpeech-100 transcripts; sentences from our  
 1230 train-validation-test splits are excluded. We sum-  
 1231 marize this approach, denoted EMG2TEXT with

<sup>5</sup>We use the WFST decoding implementation provided by ICEFALL ([github.com/k2-fsa/icefall](https://github.com/k2-fsa/icefall)).

1232

direct EMG2SPEECH in table 6.

Table 6: PER and WER for EMG2TEXT and EMG2SPEECH.

TRAINING METHOD	PER (% ↓)	WER (% ↓)
EMG2SPEECH	42.79	61.02
EMG2TEXT	38.19	51.17

1233

1234

1235

1236

1237

1238

1239

1240

1241

Overall, EMG2TEXT achieves lower PER and WER than direct EMG2SPEECH. However, direct EMG-to-speech generation remains important for neural prostheses because it can enable a more fluid, natural interaction (e.g., without requiring an explicit intermediate text interface). We therefore view improving direct EMG-to-speech as an important direction for future work.