Non-invasive electromyographic speech neuroprosthesis: a geometric perspective

Anonymous Authors¹

Abstract

In this article, we present a high-bandwidth egocentric neuromuscular speech interface for translating silently voiced speech articulations into text and audio. Specifically, we collect electromyogram (EMG) signals from multiple articulatory 015 sites on the face and neck as individuals articulate speech in an alaryngeal manner to perform EMG-to-text or EMG-to-audio translation. Such 018 an interface is useful for restoring audible speech 019 in individuals who have lost the ability to speak in-020 telligibly due to laryngectomy, neuromuscular disease, stroke, or trauma-induced damage (e.g., radiotherapy toxicity) to speech articulators. Previous works have focused on training text or speech synthesis models using EMG collected during 025 audible speech articulations or by transferring audio targets from EMG collected during audi-027 ble articulation to EMG collected during silent 028 articulation. However, such paradigms are not 029 suited for individuals who have already lost the 030 ability to *audibly* articulate speech. We are the first to present an alignment-free EMG-to-text and EMG-to-audio conversion using only EMG collected during silently articulated speech in an 034 open-sourced manner. On a limited vocabulary 035 corpora, our approach achieves almost $2.4 \times$ improvement in word error rate with a model that is $25 \times$ smaller by leveraging the inherent geometry of EMG. 039

042 **1. Introduction**

041

051

044Electromyogram (EMG) signals gathered from the orofa-
cial neuromuscular system during the silent articulation of
speech in an alaryngeal manner can be synthesized into per-
sonalized audible speech, potentially enabling individuals

without vocal function to communicate naturally. Furthermore, such systems could seamlessly interface with virtual environments where audible communication might disturb others (e.g., multiplayer games) or facilitate telephonic conversations in noisy environments. A key enabler of such advancements is the rich information encoded in EMG signals recorded from multiple spatially separated locations, which capture muscle activation patterns across different muscles. This richness allows for the decoding of subtle and intricate details, such as nuanced speech articulations, likely with higher bandwidth and lower latency compared to exocentric or allocentric modalities, such as video-based lip-to-speech synthesis. By leveraging this information, EMG-based systems offer a promising foundation for natural and efficient communication across diverse applications.

In this article, we present EMG-to-language translation models with a focus on data *geometry*. We show that EMG-tolanguage translation can be cast as a graph-connectivity learning problem and provide a *single*-layer recurrent architecture with connectionist temporal classification (CTC) loss (Graves et al., 2006) on the manifold of symmetric positive definite (SPD) matrices. Our alignment free translation method is similar to paradigms proposed for invasive speech brain-computer interfaces described by Willett et al. and Metzger et al. While invasive methods are viable for individuals with anarthria or amyotrophic lateral sclerosis, our EMG based non-invasive speech prosthesis is appropriate for individuals who have undergone laryngectomy or experience dysarthria or dysphonia.

On a limited English dataset with a vocabulary of 67 words, we demonstrate that our model achieves a decoding accuracy of 88% on word transcriptions. On a larger, general English language corpus, we achieve a phoneme error rate (PER) of 56%, as measured using Levenshtein distance (between the original and constructed phoneme sequences). Additionally, we show that the model can be trained with minimal data, achieving good performance even when tested on a dataset nearly 5 times larger than the training set, where sentences are spelled out using NATO phonetic codes. This capability is crucial, as collecting large-scale datasets for such systems is often challenging. These results highlight the potential for the practical deployment of such interfaces at scale.

 ¹Anonymous Institution, Anonymous City, Anonymous Region,
 Anonymous Country. Correspondence to: Anonymous Author
 <anon.email@domain.com>.

<sup>Proceedings of the 42nd International Conference on Machine
Learning, Vancouver, Canada. PMLR 267, 2025. Copyright 2025
by the author(s).</sup>

055 **2. Prior work**

The current benchmark in silent speech interfaces is estab-057 lished by Gaddy & Klein; Gaddy & Klein. Using elec-058 tromyogram (EMG) signals collected during silently artic-059 ulated speech (E_S) and *audibly* articulated speech (E_A) , 060 along with corresponding audio signals (A), they develop 061 a recurrent neural transduction model to map time-aligned 062 features of E_A or E_S with A. In their baseline model, joint 063 representations between E_A and A are learned during train-064 ing, and the model is tested on E_S . To improve performance, 065 a refined model aligns E_S with E_A and subsequently uses 066 the aligned features to learn joint representations with A. 067 The methods described above have significant shortcomings 068 that limit their practicality for real-world deployment. They 069 are, (1) the unavailability of good quality E_A and A in in-070 dividuals who have lost vocal and articulatory functions, (2) the need for a 2x sized training corpus for learning xrepresentations (both E_A and E_S), and (3) the requirement of aligned features, which are computationally expensive 074 and time-consuming to obtain, making near real-time imple-075 mentation challenging. We overcome the above challenges 076 by training a model with only E_S and corresponding phonetic transcription without any alignments, using CTC loss. 078 Besides, unlike Gaddy & Klein; Gaddy & Klein, we demon-079 strate the efficacy of our models on multiple subjects.

081 Another notable approach is presented by Gowda et al., who 082 demonstrate that, unlike images and audio - which are func-083 tions sampled on Euclidean grids - EMG signals are defined 084 by a set of orthogonal axes, with the manifold of SPD ma-085 trices as their natural embedding space. We build upon the 086 methods described by Gowda et al. in our analysis and 087 introduce the following key improvements: (1) we train a re-088 current model for EMG-to-phoneme sequence-to-sequence 089 generation, as opposed to the classification models proposed 090 by Gowda et al., (2) we operate in the sparse graph spectral 091 domain, effectively circumventing bottlenecks associated 092 with repeated eigenvalue computation in neural networks, 093 which, due to their iterative nature, often have limited paral-094 lelization capabilities on GPUs, and (3) demonstrate EMG-095 to-language conversion on continuously articulated speech 096 as opposed to individual words or phonemes. 097

A substantial body of prior work (Jou et al., Schultz & Wand, Kapur et al., Meltzner et al., Toth et al., Janke & Diener, and Diener et al.) has laid the groundwork for the development of silent speech interfaces. While these studies have been instrumental in shaping the field, they place less emphasis on understanding the *data structure* and the implementation of parameter and data-efficient approaches.

In the following sections, ① we explain the inherent non-Euclidean data structure of EMG signals, ② quantify the signal distribution shift across individuals, and ③ demonstrate that high fidelity phoneme-by-phoneme translation of EMG-to-language is possible using only E_S without E_A and A.

3. Methods

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

To model the mapping from E_S to L, we employ a sequenceto-sequence model trained using CTC loss. This approach allows us to train the model with *unaligned* pairs of E_S and L, eliminating the need for precise alignment between the input signals and their corresponding phoneme sequences. During testing, a sample of E_S not in the training set outputs probabilities over all possible phonemes (40 of them in our case) at every time step, and we construct L using beam search. L is then converted to personalized audio A using few-shot learning (Choi et al., 2021), which requires as little as a single audio clip from the individual (an audio clip of about 3-5 minutes, not necessarily containing the same phonemic content as L, recorded before their clinical condition). By leveraging this sample, we generate audio Athat captures both the predicted linguistic content and the speaker's unique vocal characteristics (we elaborate on this topic in appendix G).

3.1. EMG data representation

Gowda et al. demonstrate that the manifold of SPD matrices serves as an effective embedding space for EMG signals, enabling the natural distinction of different orofacial movements associated with speech articulation and all English phonemes using raw signals. We make significant improvements on their methods to perform phoneme-by-phoneme decoding as opposed to classification paradigms and demonstrate our methods on continuously articulated speech in the English language as opposed to discrete word or phoneme articulations.

We construct a complete graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}(\tau))$ to represent the functional connectivity of EMG signals, where $\mathcal{E}(\tau)$ denotes the set of edges over a time window $\tau = [t_{\text{START}}, t_{\text{END}}]$ (Gowda et al., 2024). The edge weight between two nodes v_1 and $v_2 \in \mathcal{V}$ in a time window is defined as $e_{12} = e_{21} = \mathbf{f}_{v_1}^T \mathbf{f}_{v_2}$, which corresponds to the covari110 ance of the signals at those nodes during the time interval. 111 Consequently, the edge (adjacency) matrix $\mathcal{E}(\tau)$ is sym-112 metric and positive semi-definite. Following Gowda et al., 113 we convert semi-definite adjacency matrices into definite 114 ones by adding a shrinkage estimator. We then model these 115 symmetric positive definite (SPD) matrices using the Rie-116 mannian geometry approach via Cholesky decomposition, 117 as described by Lin.

118 For any adjacency matrix \mathcal{E} , we can express it as \mathcal{E} = 119 $U\Sigma U^T$, where U is the matrix of eigenvectors, and Σ is a 120 diagonal matrix containing the corresponding eigenvalues. 121 However, instead of calculating U for each \mathcal{E} at every time-122 step τ , we fix an approximate common eigenbasis Q derived 123 from the Fréchet mean \mathcal{F} (Lin, 2019) of all adjacency matri-124 ces (at different time points) in the training set. Specifically, 125 we compute \mathcal{F} as the geometric mean of all \mathcal{E} , and decom-126 pose it as $\mathcal{F} = Q\Lambda Q^T$, where Q contains the eigenvectors 127 of \mathcal{F} , and Λ is a diagonal matrix of its eigenvalues. 128

129 Using this fixed eigenbasis Q, any adjacency matrix \mathcal{E} can 130 be approximately diagonalized as $Q^T \mathcal{E} Q$, yielding a sparse 131 matrix σ . Gowda et al. show that such a matrix Q can be 132 learned using neural networks constrained on the Stiefel 133 manifold (Huang & Van Gool, 2017) and that such a Q134 is different for different individuals. However, neural net-135 works constrained on the Stiefel manifold require perform-136 ing repeated eigendecomposition operations, which have 137 limited parallelization capability and lead to unstable gra-138 dients while using CTC loss. Therefore, we simply derive 139 Q from the Fréchet mean \mathcal{F} and use that Q to obtain sparse 140 matrices σ . In appendix D, we show that matrices σ are 141 indeed sparse by comparing the ratio of maximum value of 142 non-diagonal entries to maximum value of diagonal entries 143 of matrices before and after approximate diagonalization. 144 This formulation allows us to work in an approximate graph 145 spectral domain with a consistent orthogonal basis across all time windows τ . For our task, we compute the graph 147 spectral sequences σ for all time windows τ and use these 148 as inputs for EMG-to-language translation. We illustrate 149 these concepts in figure 1.

151 3.2. Sequence-to-sequence modeling with a single recurrent layer

We implement a single-layer gated recurrent unit (GRU) architecture for EMG-to-phoneme sequence-to-sequence modeling. The input to the GRU consists of a sequence of approximately diagonalized matrices, denoted as σ , derived over different time windows τ .

¹⁵⁹ To investigate whether recurrent models defined on the manifold provide a better representation of $\sigma(\tau)$ compared to those defined in Euclidean space, we construct three distinct GRU architectures:



Figure 1. Conceptual depiction of SPD edge matrices on a 3D convex cone manifold. Edge matrices derived from a given individual can be represented using an approximate common eigenbasis. Consider six edge matrices, corresponding to six different time windows τ_1 to τ_6 , for an individual A. These matrices are shown in *blue*. The Fréchet mean of these SPD matrices, denoted as \mathcal{F}_A , is represented in *pink*. Each edge matrix $\mathcal{E}_A(\tau_i)$ can be expressed as: $\mathcal{E}_A(\tau_i) = U_{A_i} \Sigma_{A_i} U_{A_i}^T$, where Σ_{A_i} is a diagonal matrix of eigenvalues, and U_{A_i} is an orthogonal matrix of eigenvectors ($i \in [1, 6]$). Instead of representing each edge matrix with a separate eigenvector matrix, we transform all $\mathcal{E}_A(\tau_i)$ using a *common basis* Q_A , which corresponds to the eigenvectors of the Fréchet mean \mathcal{F}_A . Specifically, we calculate $\sigma_{A_i} = Q_A^T \mathcal{E}_A(\tau_i) Q_A$. Matrices σ_{A_i} are approximately diagonal. We use these approximately diagonalized representations for EMG-to-language translation. Separately, edge matrices from another individual B, shown in green, have their Fréchet mean represented in yellow. Matrices belonging to different individuals reside in distinct neighborhoods on the manifold, and the common basis Q_A for individual A cannot approximately diagonalize edge matrices from individual B. Instead, a separate basis Q_B , derived from the eigenvectors of \mathcal{F}_B , is required for individual B. Geometrically, the manifold is only locally Euclidean and the tangent spaces for individuals A and B are distinct. That is, transformation of the space $\mathbb{R}^{|\mathcal{V}|}$ induced by EMG signals of subjects A and B are different and the approximate orthogonal eigenbasis vectors that characterize such transformations are different for different individuals. This signal distribution shift can be approximated as change of basis. An inset diagram illustrates the eigenvectors of $\mathcal{E}_A(\tau_i)$.

(1) **GRU**_A: A GRU layer defined in the Euclidean domain, following the implementation described by Chung et al. (2014),

(2) \mathbf{GRU}_B : A GRU layer formulated on the manifold of SPD matrices, as proposed by Jeong et al. (2024), and

③ **GRU**_C: A GRU layer defined on the manifold of SPD matrices, plus an implicit layer solved using neural ordinary differential equations, integrating methodologies from Jeong et al. (2024), Chen et al. (2018), and Lou et al. (2020).

164

 GRU_B and GRU_C directly accept SPD matrices, σ , as input, 165 166 whereas GRU_A processes the vectorized representations of 167 σ . At each time step, the GRU models output probability 168 distributions over 40 phonemes in the English language. The 169 models are trained using CTC loss, and during inference, 170 the most probable phoneme sequence is reconstructed using 171 beam search decoding. The end-to-end EMG-to-language 172 translation model is depicted in figure 2.



Figure 2. Illustration of multivariate EMG-to-phoneme sequence translation. Bandpass-filtered and z-normalized raw signals are converted to SPD edge matrices, $\mathcal{E}(\tau)$, over a time window τ . These edge matrices are then transformed into approximately diagonalized matrices, $\sigma(\tau)$, which are fed into a bidirectional GRU layer. At each time step (every 20 ms), the GRU outputs probability distributions *P* over 40 phonemes in the English language. During inference, the most probable phoneme sequence is reconstructed using beam search decoding.

196

206

3.3. Geometric perspective aligns well with biology

208 We study multivariate EMG signals collected at $|\mathcal{V}|$ sensor 209 nodes in different time windows τ using edge matrices $\mathcal{E}(\tau)$, 210 which capture the relationship between every pair of nodes 211 in $|\mathcal{V}|$.

This can be understood as studying the transformation of the space $\mathbb{R}^{|\mathcal{V}|}$ given by the transformation matrix $\mathcal{E}(\tau)$. Such a transformation can equivalently be expressed in a coordinate system where the eigenvectors serve as the basis vectors. This change of basis is described by $U^T \mathcal{E}(\tau) U = \Sigma(\tau)$, where $\Sigma(\tau)$ is a diagonal matrix. In this eigenbasis coordinate system, transformation of space induced by EMG signals can be interpreted as a linear combination of the columns of U, with the diagonal values of Σ acting as coefficients. By fixing an approximate eigenbasis Q, we obtain an approximately diagonal matrix $\sigma(\tau)$ and an approximate linear combination. This formulation aligns well with the biological process underlying EMG signal generation, which involves a purely additive combination of muscle activations which contrasts with processes such as speech production, which can be modeled as the application of a time-varying filter to a time-varying source signal (Sivakumar et al.). For a given individual, we fix the approximate eigenbasis vectors and focus on analyzing the approximate eigenvalues $\sigma(\tau)$. These matrices can then be studied using a single layer recurrent neural network. EMG signals from different individuals induce different transformations of $\mathbb{R}^{|\mathcal{V}|}$ and have different eigenbasis vectors. The signal distribution shift across different individuals can thus be interpreted as change of basis in $\mathbb{R}^{|\mathcal{V}|}$. It should be noted that while a shallow single layer network is sufficient to learn multivariate EMG-phoneme translation using sparse matrices σ , and while such an architecture works consistently well across subjects, the weights of the recurrent networks must be fine-tuned for different individuals, as σ from different individuals correspond to different basis vectors Q.

4. Data

We evaluate our models using three datasets, referred to as Data $_{\text{SMALL-VOCAB}}$, Data $_{\text{LARGE-VOCAB}}$, and Data $_{\text{NATO-WORDS}}$, which are described below. We use Data $_{\text{SMALL-VOCAB}}$ and Data $_{\text{LARGE-VOCAB}}$ to evaluate naturally articulated speech in a silent manner. We use Data $_{\text{NATO-WORDS}}$ to demonstrate that, by using a small codeword set such as NATO codes, we can construct a generalizable language-spelling model that requires very little data for training. Additionally, Data $_{\text{NATO-WORDS}}$ is used to show that our models work consistently well across individuals. This paradigm is useful for rapid training (or fine-tuning) and deployment of speech prostheses.

4.1. Data SMALL-VOCAB

Following Gaddy & Klein, we create a limited vocabulary dataset consisting of 67 unique words. These words include weekdays, ordinal dates, months, and years. Sentences are constructed from these words in the format \langle WEEKDAY-MONTH-DATE-YEAR>. A single individual articulated 500 such sentences silently, and the resulting EMG data, E_S , is translated into output phoneme sequences. We have timestamps to demarcate the beginning and the end of words within a sentence.

We collect EMG data from 31 muscle sites at a sampling rate of 5000 Hz. Of these, 22 electrode sites are identical to those used by Gowda et al., while the remaining 9 electrodes are placed symmetrically on the opposite side of the neck. The experimental setup is same as that described by Gowda et al.

4.2. Data LARGE-VOCAB

We adapt the language corpora from Willett et al., who demonstrated a speech brain-computer interface by translating neural spikes from the motor cortex into speech. The dataset comprises an extensive English language corpus containing approximately 6,500 unique words and 11,000 sentences. Unlike Gaddy & Klein; Gaddy & Klein, we collect only E_S (excluding E_A and A) and perform E_S to-language translation without time-aligning with E_A and A. The data collection setup follows the methodology described for Data SMALL-VOCAB. This corpus includes sentences of varying lengths, with the subject articulating sentences at a normal speed, averaging 160 words per minute. Timestamps were used solely to mark the beginning and end of each sentence, with the subject clicking the mouse at the start of articulation and again upon completion (unlike Data SMALL-VOCAB, there are no timestamps to demarcate between words within a sentence).

4.3. Data NATO-WORDS

We use the dataset provided by Gowda et al.¹ Specifically, we use data from their second experiment, in which 4 individuals articulated English sentences in a spelled-out manner using NATO phonemic codes in a silent manner. For instance, the word <RAINBOW> was articulated as <ROMEO-ALFA-INDIA-NOVEMBER-BRAVO-OSCAR-WHISKEY> with phonemic transcription <R-OW-M-IY-OW|AE-L-F-AH | IH-N-D-IY-AH | N-OW-V-EH-M-B-ER | B-R-AA-V-OW | AO-S-K-ER | W-IH-S-K-IY>. Subjects articulated phonemically balanced RAINBOW and GRAND-FATHER passages in this spelled-out format. In total, 1968 NATO code articulations were recorded across both passages.

The EMG data was collected from 22 muscle sites in the neck and cheek regions at a sampling rate of 5000 Hz. We present results for Data NATO-WORDS in appendix A.

5. Results

We describe the experimental setup and results for Data _{SMALL-VOCAB} and Data _{LARGE-VOCAB}, providing a comparative analysis with previous benchmarks.

5.1. Results for Data SMALL-VOCAB

Raw EMG signals are bandpass filtered between 80 and 1000 Hz and are *z*-normalized per channel along the time

¹The dataset is available at Gowda et al. dataset.

dimension. Then, a complete time dependent graph ($\mathcal{E}(\tau)$) and $\sigma(\tau)$) is constructed using the EMG signals. We follow the same *train-validation-test* split outlined by Gaddy & Klein. All parameters are detailed below in table 1.

Table 1. Experimental setup for Data SMALL-VOCAB.

	Data _A properties
τ	50 ms (a sliding window with an
	overlapping context size of 100 ms
	and a step size of 50 ms)
$\mathcal{E}(\tau)$ and $\sigma(\tau)$	SPD matrices of dimensions 31×31
Train-validation-test split	370 - 30 - 100 sentences
Beamsearch width	Top-5

The Fréchet mean, computed from the training set, is utilized to calculate $\sigma(\tau)$ for all τ in the training, validation, and test datasets. Figure 3 illustrates the Levenshtein distances between target and predicted phoneme sequences for three GRU models with varying model sizes. To decode the articulated word(s), we identify a word or a set of words from the vocabulary corpus whose phonemic sequence best matches the predicted sequence, using Levenshtein distance as the metric.

Decoding accuracy for EMG-to-text translation is evaluated as 1 - WER and is presented in figure 4 for models of different sizes. Model size is controlled exclusively by adjusting the GRU hidden unit dimensionality, which is the *only hyperparameter* in our setup. On this limited vocabulary corpus, we achieve a WER as low as 12%, with the average Levenshtein distance between target and predicted sequences below 1. These results underscore the feasibility and practical potential of EMG-to-language translation technology.



Figure 3. Model size versus Levenshtein distance. Models are evaluated over 10 random seeds. *Lower is better.*



287 Figure 4. Decoding accuracy = 1 - WER versus model size. Models are evaluated over 10 random seeds. Higher is better.

291 Now, we compare our results with the results given by 292 Gaddy & Klein. Gaddy & Klein recorded EMG signals us-293 ing 8 electrodes at a sampling rate of 1000 Hz. To enable a 294 direct comparison with their approach, we downsample our 295 EMG signals from 5000 Hz to 1000 Hz and select a subset 296 of 8 electrodes from the original 31. The placement of these 297 electrodes is approximately aligned with those specified by Gaddy & Klein to ensure consistency in the experimental 299 setup. 300

We compare our results with their baseline approach, where 301 models were trained for E_A -to-A translation and evaluated 302 on E_S -to-A translation. In contrast, our approach empha-303 sizes direct E_S -to-L translation. Additionally, their im-304 proved framework incorporates both E_A and E_S signals, 305 relying on time alignment between them. However, they do 306 not propose a paradigm that independently translates E_S 307 without leveraging E_A or A. In this case, $\mathcal{E}(\tau)$ and $\sigma(\tau)$ 308 are 8×8 matrices. The rest of the training paradigm is same 309 as in table 1. We provide the comparison in table 2. Our 310 approach achieves almost $2.4 \times$ improvement in WER with 311 a model that is $25 \times$ smaller. 312

314 Table 2. Comparison with Gaddy & Klein. Our approach achieves 315 almost $2.4 \times$ improvement in WER with a model that is $25 \times$ 316 smaller. WER is averaged over 10 random seeds. Results are for 317 Data_{SMALL-VOCAB} using 8 electrodes with signals downsampled to 318 1000 Hz. 319

320		
321	Ours	Baseline of Gaddy &
322		Klein
323	WER - 27%, using	WER - 64%
324	GRU_B	
325	Model size - about	Model size - about 40
326	1.4 million	million
327		

313

289 290

3

3 32

329

5.2. Results for Data LARGE-VOCAB

As in Data SMALL-VOCAB, we filter, z-normalize, and construct $\sigma(\tau)$. The properties of the dataset are detailed in table 3. Sentences in the validation and test sets do not occur in the training set. On this general English language corpus consisting of approximately 6500 words, spoken at an average rate of 160 words per minute, we perform EMG-to-phoneme sequence translation and measure the Levenshtein distances between the target and predicted phoneme sequences. The phoneme error rates (PER) are presented in table 4. Transcription examples are given in table 5.

Table 3. Experimental setup for Data LARGE-VOCAB.

	Data LARGE-VOCAB properties
au	20 ms (a sliding window with an
	overlapping context size of 50 ms
	and a step size of 20 ms)
$\mathcal{E}(\tau)$ and $\sigma(\tau)$	SPD matrices of dimensions 31×31
Train-validation-test split	8000 - 1000 - 1970 sentences
Beamsearch width	Top-5

Table 4. We achieve a PER of 56% for speech articulated at 160 words per minute using 31 electrodes. Average phoneme sequence length of sentences is 24.5, and the chance decoding accuracy of a sequence (that is, chance 1 - PER) is $\left(\frac{1}{40}\right)^{24.5}$. While the results are modest compared to high density invasive speech brain-computer interfaces, we showcase significant potential of a non-invasive method. Willett et al. report a PER of 21% using 128 intracortical arrays at a slower speech rate of 62 words per minute. Metzger et al. report a PER of 30% on a smaller corpora of 1024 words articulated at a slower rate of 78 words per minute using 253 ECoG electrodes. In future work, we would like to verify if higher density EMG and more training data can lead to better PER.

Ours - 31 electrodes
PER - 56%, using GRU_A
Model size - about 4.4 mil-
lion

6. Observations and discussions

(1) From Data SMALL-VOCAB and Data LARGE-VOCAB and their corresponding results, we observe that continuously articulated silent speech - where subjects naturally articulate sentences in English but inaudibly - can be translated into phonemic sequences at a fine-scale resolution of 50 ms or 20 ms. This resolution is comparable to state-of-the-art (SOTA)

automatic speech recognition (ASR) models, such as those described by Baevski et al. and Hsu et al., which operate at a 20 ms resolution. These findings highlight the potential 333 for real-time EMG-to-language translation, akin to audio-334 to-audio (language translation) and audio-to-text translation. 335 Furthermore, achieving a word transcription accuracy of approximately 88% on limited-vocabulary corpora and a PER 337 of 56% on open-vocabulary corpora demonstrates that high-338 fidelity translation is achievable, reinforcing the viability of 339 this approach.

340 (2) Additionally, results from Data NATO-WORDS (WER of 59%) 341 averaged across all 4 subjects) indicate that by leverag-342 ing NATO phonetic codes, we can establish a generaliz-343 able EMG-to-language spelling paradigm. Although this approach does not replicate natural speech, it enables a 345 practical mode of limited communication for individuals 346 who have lost speech articulation capabilities. Notably, 347 this paradigm is efficient, requiring only a small corpus for training - our model trained on just 10 minutes of data, 349 demonstrates robust generalization on a much larger test set. 350

351 (3) The key contribution of this article lies in the devel-352 opment of efficient architectures for multivariate EMG-to-353 phoneme sequence translation. We show that EMG signals 354 can be approximately decomposed into linear combinations 355 of a set of orthogonal axes, represented by the matrices 356 $\sigma(\tau)$. This decomposition enables the analysis of time-357 varying graph edges in a sparse graph spectral domain using 358 a single recurrent layer. Notably, our model relies on only 359 one hyperparameter - the dimension of the hidden unit in 360 the GRU. Across different datasets and subjects, the models 361 exhibit consistent and predictable behavior with respect to the hidden unit dimension. Specifically, the decoding ac-363 curacy of GRU_A and GRU_B improves as the hidden unit dimension increases, eventually plateauing, while GRU_C 365 demonstrates a peak in performance before diminishing 366 (figures 4 and 6). Also, GRU_C outperforms other GRU 367 models for Data NATO-WORDS. It also outperforms other GRU models for Data SMALL-VOCAB at smaller model sizes. This 368 369 demonstrates that modeling dynamics of EMG signals using 370 neural ODEs is beneficial and allows for better abstraction 371 of the data. Importantly, although different datasets and 372 individuals are characterized by distinct orthogonal basis 373 vectors, the same model architecture can be applied across 374 individuals without the need for extensive hyperparameter 375 tuning.

We achieve a word error rate (WER) of approximately
We achieve a word error rate (WER) of approximately
12% on a 67-word vocabulary, not too far from the 9.1%
WER reported by Willett et al. on a 50-word vocabulary.
Notably, our results are achieved using only 31 non-invasive
electrodes, in contrast to the 128 intracortical electrodes
employed by Willett et al. On Data LARGE-VOCAB, we achieve
a phoneme error rate (PER) of 56% for speech articulated

384

at an average rate of 160 words per minute, whereas Willett et al. report a PER of 21% using 128 intracortical arrays at a slower speech rate of 62 words per minute. In future work, we would like to verify if higher density EMG and more training data can lead to better PER. These findings demonstrate the feasibility of a non-invasive approach for translating silent speech into language. While Willett et al. and Metzger et al. showcase brain-computer speech interfaces for individuals with anarthria or amyotrophic lateral sclerosis, our method provides a viable alternative for individuals who have undergone laryngectomy or experience dysarthria or dysphonia, where invasive recordings may not be a practical solution. We highlight the significant potential of non-invasive techniques for broad clinical applicability.

(5) Défossez et al. demonstrate methods for decoding speech perception from non-invasive neural recordings using magnetoencephalography (MEG) and electroencephalography (EEG). They show that *listened* speech segments can be predicted from MEG with an accuracy of 41%. However, such interfaces are not useful for initiating communication. We go beyond these models to demonstrate that, using noninvasive EMG signals, we can decode speech articulation at the phonemic level with a higher accuracy of 44% on a large English language corpus.

(6) We envision EMG-based non-invasive neuroprostheses having a user-friendly form factor that is easy to don and doff. However, even minor variations in electrode placement introduce a covariate signal shift, which can be mathematically represented as a change of basis (Gowda et al., 2024). Additional factors, such as variations in subcutaneous fat and changes in neural drive characteristics, contribute to covariate signal drift over time. To address these challenges, modeling EMG signals using SPD covariance matrices proves advantageous. Our models show consistent performance across subjects, as demonstrated here, and outperform Euclidean-space models (Gaddy & Klein, Gaddy & Klein) in terms of both decoding accuracy and model parameter efficiency. Moreover, considering the idiosyncrasies of individuals, the difficulty of collecting large-scale data, and the need for frequent fine-tuning due to circumstantial variations, a streamlined approach is crucial. A simple model leveraging a single GRU layer, as presented here, offers an effective solution for adaptability.

7. Conclusion

We present an efficient data representation for orofacial EMG signals and demonstrate that our approach enables effective EMG-to-language translation. Our method outperforms previous benchmarks on limited-vocabulary corpora, showcasing its potential for practical applications. Notably, we demonstrate the ability to translate EMG collected during *silently* voiced speech (E_S) to language without requiring

385 386 387 388 389 390	corresponding audio (A) and EMG collected during <i>audibly</i> voiced speech (E_A), marking a significant advancement in the translation paradigm and paving the way for real-world deployment of such devices. By providing open-source data and code, this work lays a solid foundation for the development of efficient neuromuscular speech prostheses.
 391 392 393 394 205 	In future work, we plan to augment our methods with lan- guage models and test their applicability for individuals with clinical etiologies that affect voicing and articulator movement in real time.
395 396	
397	
398	
399	
400	
401	
402 403 404	Table 5. Examples of EMG-to-phoneme sequence translations. We do translations using EMG collected during <i>silent</i> articulations (E_S) with CTC loss without making use of corresponding time aligned <i>audio</i> (A) and EMG collected during <i>audible</i> articulation (E_A) . Ground truth sentences with corresponding timestamps. Ground truth phonemic transcriptions. Decoded phonemic transcriptions.
405	3 transcribed sentences in Data SMALL-VOCAR
400	T-START <wednesday>_{T-END} T-START <july>_{T-END} T-START <twenty sixth="">_{T-END} T-START <nineteen seven="" sixty="">_{T-END} W-EH-N-Z-D-IY SPACE J-UW-L-AY SPACE T-W-EH-N-T-IY-S-IH-K-S-TH SPACE N-AY-N-T-IY-N-S-IH-K-S-T-IY-S-EH-V-AH-N</nineteen></twenty></july></wednesday>
408	W-AH-N-Z-D-IY SPACE J-UW-L-AY SPACE T-W-EH-N-T-IY-S-IH-K-S-TH SPACE N-AY-N-T-IY-N-S-IH-K-S-T-IY-S-EH-V-AH-N
410	T-START <thursday>_{T-END T-START} <october>_{T-END T-START} <twenty ninth="">_{T-END T-START} <two nine="" thousand="">_{T-END} TH-FR-7-D-FY SPACE A A-K-T-OW-R-FR SPACE T-W-FH-N-T-IV-N-AY-N-TH SPACE T-UW-TH-AW-7-AH-N-D-N-AY-N</two></twenty></october></thursday>
411	TH-ER-Z-D-EY SPACE AA-K-T-OW-B-ER SPACE T-W-EH-N-T-IY-N-AY-N-TH SPACE T-UW-TH-AW-Z-AH-N-D-T-N-AY-N
41Z //13	T-START <tuesday>T-END T-START <december>T-END T-START <fifth>T-END T-START <nineteen eight="" seventy="">T-END</nineteen></fifth></december></tuesday>
414	T-UW-Z-D-IY SPACE D-IH-S-EH-M-B-ER SPACE F-IH-F-TH SPACE N-AY-N-T-IY-N-S-EH-V-AH-N-T-IY-EY-T
415	T-UW-Z-D-IY SPACE D-IH-S-EH-M-B-ER SPACE F-IH-F-TH SPACE N-AY-N-T-IY-N-S-EH-V-AH-N-T-IY-AY-N-T
416	Top-3 (best) transcribed sentences in Data LARGE-VOCAB
417	T-START <its fun="" kind="" of="">_{T-END}</its>
418	IH-T-S SPACE K-AY-N-D SPACE AH-V SPACE F-AH-N
419	IH-T-S space K-AY-N-D space F-AH-N
420	T-START <probably seventies="">T-END</probably>
421	P-R-AA-B-AH-B-L-IY SPACE S-EH-V-AH-N-T-IY-Z
422	P-R-AH-B-L-IY SPACE S-EH-V-AH-N-T-IY
423	_{T-START} <is it="" legend="">_{t-end}</is>
424	IH-Z SPACE IH-T SPACE L-EH-JH-AH-N-D
425	IH-T SPACE IH-T SPACE S-EH-JH-AH-N
426	Bottom-3 (worst) transcribed sentences in Data LARGE-VOCAB
427	_{t-start} <member american="" meteorological="" of="" society="" the="">_{t-end}</member>
428	M-EH-M-B-ER SPACE AH-V SPACE DH-AH SPACE AH-M-EH-R-AH-K-AH-N SPACE
429	M-IY-T-IY-AO-R-AH-L-AA-JH-IH-K-AH-L SPACE S-AH-S-AY-AH-T-IY
430	DH-AH-M-AH SPACE F-AH-B-AE-I-AH SPACE UW SPACE K-L SPACE S-AH SPACE I-IY SPACE D-IH
431	T-START $<$ PICTURES AND PROJECTS THAT TOUCAN MAKE TOURSELF $>_{\text{T-END}}$
432 422	P-IH-K-CH-EK-Z SPACE AH-N-D SPACE P-K-AA-JH-EH-K-I-S SPACE DH-AE-I SPACE Y-UW SPACE
433 131	K-AE-N SPACE M-EI-K SPACE I-EK-S-EH-L-F
434 435	THE SPACE I TENTE SPACE ITENTS AND SPACE FOR AN AND SPACE DOTIONS AT THE FUND OF THE VEAD
436	T-START STILL DATE OF THE DAY OF THE LEAV T-END HH-IY SPACE IH-K-S-P-FH-K-T-AH-D SPACE K-AH-N-K-I -IW-7H-AH-N-7 SPACE AF-T SPACE
437	DH-AH SPACE FH-N-D SPACE AH-V SPACE DH-AH SPACE Y-IH-R
438	G-EH-P-IH SPACE AH-K-L-UW-ZH-AH SPACE AH-N SPACE AY-D SPACE AH SPACE TH-Y
439	C EM T IN SUCCESSING REE ON EM INTEGRACE FAIT LY SURCE FAIT E SURCE FAIT SURCE FAIT I

440 Impact statement

This article provides data and methods for developing noninvasive EMG based neuroprosthesis. Such devices have the potential to restore natural communication in individuals who have lost the ability to speak intelligibly due to causes such as neuromuscular disease, stroke, trauma, and head/neck cancer surgery (e.g. laryngectomy) or treatment (e.g. radiotherapy toxicity to the speech articulators).

IRB was approved for human subject research. IRB details and its terms and conditions will be made available if the manuscript is accepted.

We are unable to publish the data and codes in an anonymized format. Data and codes will be made publicly available if the manuscript is accepted.

References

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

- Arsigny, V., Fillard, P., Pennec, X., and Ayache, N. Geometric means in a novel vector space structure on symmetric positive-definite matrices. *SIAM Journal on Matrix Analysis and Applications*, 29(1):328–347, 2007. doi: 10.1137/050637996. URL https://doi.org/10.1137/050637996.
- Baevski, A., Zhou, Y., Mohamed, A., and Auli, M. wav2vec 2.0: A framework for self-supervised learning of speech representations. *Advances in neural information processing systems*, 33:12449–12460, 2020.
- Barachant, A., Bonnet, S., Congedo, M., and Jutten, C. Multiclass brain–computer interface classification by riemannian geometry. *IEEE Transactions on Biomedical Engineering*, 59(4):920–928, 2011.
- Barachant, A., Bonnet, S., Congedo, M., and Jutten, C. Classification of covariance matrices using a riemannianbased kernel for bci applications. *Neurocomput.*, 112: 172–178, July 2013. ISSN 0925-2312. doi: 10.1016/j. neucom.2012.12.039. URL https://doi.org/10. 1016/j.neucom.2012.12.039.
- Chen, R. T., Rubanova, Y., Bettencourt, J., and Duvenaud, D. K. Neural ordinary differential equations. *Advances in neural information processing systems*, 31, 2018.
- Choi, H.-S., Lee, J., Kim, W., Lee, J., Heo, H., and Lee, K. Neural analysis and synthesis: Reconstructing speech from self-supervised representations. *Advances in Neural Information Processing Systems*, 34:16251–16265, 2021.
- Chung, J., Gulcehre, C., Cho, K., and Bengio, Y. Empirical
 evaluation of gated recurrent neural networks on sequence
 modeling. *arXiv preprint arXiv:1412.3555*, 2014.

- Défossez, A., Caucheteux, C., Rapin, J., Kabeli, O., and King, J.-R. Decoding speech perception from noninvasive brain recordings. *Nature Machine Intelligence*, 5(10):1097–1107, 2023.
- Diener, L., Felsch, G., Angrick, M., and Schultz, T. Sessionindependent array-based emg-to-speech conversion using convolutional neural networks. In *Speech Communication; 13th ITG-Symposium*, pp. 1–5, 2018.
- Gaddy, D. and Klein, D. Digital voicing of silent speech. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 5521–5530, 2020.
- Gaddy, D. and Klein, D. An improved model for voicing silent speech. In *Proceedings of the 59th Annual Meeting* of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers), pp. 175–181, 2021.
- Gowda, H. T. and Miller, L. M. Topology of surface electromyogram signals: hand gesture decoding on riemannian manifolds. *Journal of Neural Engineering*.
- Gowda, H. T., McNaughton, Z. D., and Miller, L. M. Geometry of orofacial neuromuscular signals: speech articulation decoding using surface electromyography. arXiv preprint arXiv:2411.02591, 2024.
- Graves, A., Fernández, S., Gomez, F., and Schmidhuber, J. Connectionist temporal classification: labelling unsegmented sequence data with recurrent neural networks. In *Proceedings of the 23rd international conference on Machine learning*, pp. 369–376, 2006.
- Hsu, W.-N., Bolte, B., Tsai, Y.-H. H., Lakhotia, K., Salakhutdinov, R., and Mohamed, A. Hubert: Selfsupervised speech representation learning by masked prediction of hidden units. *IEEE/ACM transactions on audio*, *speech, and language processing*, 29:3451–3460, 2021.
- Huang, Z. and Van Gool, L. A riemannian network for spd matrix learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31, 2017.
- Janke, M. and Diener, L. Emg-to-speech: Direct generation of speech from facial electromyographic signals. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 25(12):2375–2385, 2017. doi: 10.1109/TASLP.2017.2738568.
- Jeong, S., Ko, W., Mulyadi, A. W., and Suk, H.-I. Deep Efficient Continuous Manifold Learning for Time Series Modeling . *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 46(01):171–184, January 2024. ISSN 1939-3539. doi: 10.1109/TPAMI.2023.3320125. URL https://doi.ieeecomputersociety. org/10.1109/TPAMI.2023.3320125.

- Jou, S.-C., Schultz, T., Walliczek, M., Kraft, F., and Waibel,
 A. Towards continuous speech recognition using surface
 electromyography. In *Ninth International Conference on Spoken Language Processing*, 2006.
- Kapur, A., Sarawgi, U., Wadkins, E., Wu, M., Hollenstein,
 N., and Maes, P. Non-invasive silent speech recognition in
 multiple sclerosis with dysphonia. In *Machine Learning for Health Workshop*, pp. 25–38. PMLR, 2020.
- Lin, Z. Riemannian geometry of symmetric positive definite matrices via cholesky decomposition. *SIAM Journal on Matrix Analysis and Applications*, 40(4):1353–1370, 2019.
- Lou, A., Lim, D., Katsman, I., Huang, L., Jiang, Q., Lim,
 S. N., and De Sa, C. M. Neural manifold ordinary differential equations. *Advances in Neural Information Processing Systems*, 33:17548–17558, 2020.
- Meltzner, G. S., Heaton, J. T., Deng, Y., De Luca, G., Roy,
 S. H., and Kline, J. C. Development of semg sensors
 and algorithms for silent speech recognition. *Journal of neural engineering*, 15(4):046031, 2018.
- Metzger, S. L., Littlejohn, K. T., Silva, A. B., Moses, D. A.,
 Seaton, M. P., Wang, R., Dougherty, M. E., Liu, J. R., Wu,
 P., Berger, M. A., et al. A high-performance neuroprosthesis for speech decoding and avatar control. *Nature*,
 620(7976):1037–1046, 2023.
- Panayotov, V., Chen, G., Povey, D., and Khudanpur, S.
 Librispeech: An asr corpus based on public domain audio books. In 2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP), pp. 5206–5210, 2015. doi: 10.1109/ICASSP.2015.7178964.
- Sabbagh, D., Ablin, P., Varoquaux, G., Gramfort, A., and
 Engemann, D. A. *Manifold-regression to predict from MEG/EEG brain signals without source modeling*. Curran
 Associates Inc., Red Hook, NY, USA, 2019.
- Schultz, T. and Wand, M. Modeling coarticulation in emg based continuous speech recognition. *Speech Communi- cation*, 52(4):341–353, 2010.
- Sivakumar, V., Seely, J., Du, A., Bittner, S. R., Berenzweig, A., Bolarinwa, A., Gramfort, A., and Mandel, M. I. emg2qwerty: A large dataset with baselines for touch typing using surface electromyography. In *The Thirty-eight Conference on Neural Information Processing Systems Datasets and Benchmarks Track.*
- Toth, A. R., Wand, M., and Schultz, T. Synthesizing speech from electromyography using voice transformation techniques. In *Interspeech 2009*, pp. 652–655, 2009. doi: 10.21437/Interspeech.2009-229.

- Veaux, C., Yamagishi, J., and King, S. The voice bank corpus: Design, collection and data analysis of a large regional accent speech database. In 2013 International Conference Oriental COCOSDA held jointly with 2013 Conference on Asian Spoken Language Research and Evaluation (O-COCOSDA/CASLRE), pp. 1–4, 2013. doi: 10.1109/ICSDA.2013.6709856.
- Willett, F. R., Kunz, E. M., Fan, C., Avansino, D. T., Wilson, G. H., Choi, E. Y., Kamdar, F., Glasser, M. F., Hochberg, L. R., Druckmann, S., et al. A high-performance speech neuroprosthesis. *Nature*, 620(7976):1031–1036, 2023.

550 A. Results for Data NATO-WORDS

Raw EMG signals are bandpass filtered between 80 and
1000 Hertz and are *z*-normalized per channel along the time
dimension. Then, a complete time dependent graph is constructed using the EMG signals. We follow the same *train- validation-test* split outlined by Gowda et al. All parameters
are detailed in table 6.

558 Like before, the Fréchet mean, computed from the training 559 set, is utilized to calculate $\sigma(\tau)$ for all τ in the training, vali-560 dation, and test datasets. Figure 5 illustrates the Levenshtein 561 distances between target and predicted phoneme sequences 562 for three GRU models with varying model sizes across all 563 4 individuals. To decode the articulated NATO phonetic 564 code, we identify a code from the corpus of 26 codes whose 565 phonetic sequence best matches the predicted sequence, us-566 ing Levenshtein distance as the metric. Decoding accuracy 567 for EMG-to-text translation is evaluated as 1 - WER and 568 is presented in figure 6 for models of different sizes across 569 all four subjects. For each subject, the best decoding accu-570 racy across all three GRU models and all model sizes are 571 summarized in table 7. 572

Table 6. Experimental setup for Data NATO-WORDS.

573

574

593 594

575		
576		Data NATO-WORDS properties
577	au	30 ms (a sliding window with an
578		overlapping context size of 150 ms
579		and a step size of 30 ms).
580	$\mathcal{E}(\tau)$ and $\sigma(\tau)$	SPD matrices of dimensions $22 \times$
581		22.
582	Train set	416 NATO alphabet articulations
583		(26 words, each repeated 16 times)
584	Validation set	104 NATO alphabet articulations (26
585		words, each repeated 4 times)
586	Test set	1968 NATO alphabet articulations
587		(entire GRANDFATHER and RAIN-
588		BOW passages articulated in a
589		spelled-out manner).
590	Beamsearch	Top-5
591	width	1
592		

595 We compare our results obtained by constructing the most probable phoneme sequences using beam search over the 596 output probability distributions at every time step to that of 597 classification models presented by Gowda et al. in table 7. 598 599 We see a slight decrease in decoding accuracy. This might be due to the fact that these are single word articulations 600 and classification models have the context of the articulation 601 duration of the entire word as opposed to 150 ms context size 602 in phoneme-by-phoneme sequence-to-sequence modeling. 603 604

Table 7. Best decoding accuracy across all GRU models and model sizes of all four subjects (calculated by averaging over 10 random seeds). Random chance accuracy is only **3.85%**.

Subject	Decoding accuracy (1 - WER) via phoneme-by-phoneme reconstruction	Gowda et al. classification model accuracy
1	44.29%	51.34%
2	44.97%	42.89%
3	29.58%	37.21%
4	43.59%	42.79%
Average	40.61 %	43.56%



Figure 5. Average Levenshtein distance between target and predicted phoneme sequences of all four subjects (subject 1 to subject 4, starting from top left in a clockwise manner). Models are evaluated over 10 random seeds. *Lower is better*.



Figure 6. Average decoding accuracy = 1 - WER of all four subjects (subject 1 to subject 4, starting from top left in a clockwise manner). Models are evaluated over 10 random seeds. *Higher is better.*

B. Training paradigms

661 We trained our models on the training set and validated 662 them on the validation set. For testing on the test set, we 663 selected the model weights corresponding to the epoch with 664 the minimum validation loss, provided the losses in the im-665 mediately preceding and succeeding epochs did not exceed 666 more than 20% of the minimum loss. All GRU models 667 were trained for 100 epochs on Data SMALL-VOCAB and 500 668 epochs on Data LARGE-VOCAB. For Data NATO-WORDS, GRUA 669 was trained for 250 epochs, while GRU_B and GRU_C were 670 each trained for 150 epochs. Model training is completed in 671 just a few minutes, whether using a GPU or CPU. 672

673 For WER calculation, we compute the Levenshtein dis-674 tance between the predicted phoneme sequence and all 675 entries in the word corpora (or combinations of them). 676 When multiple entries share the lowest distance with the 677 predicted sequence and the ground truth is among them, 678 we classify it as a wrong prediction. For instance, the 679 predicted sequence <T-TH-UW-AH-Z-D-EY> could be de-680 coded as either <TUESDAY> or <THURSDAY>, whose 681 phonetic transcriptions are <T-UW-Z-D-IY> and <TH-ER-682 Z-D-EY>, respectively. In this case, the ground truth text 683 prompt is <THURSDAY>. However, both <TUESDAY> 684 and <THURSDAY> yield the same minimum Levenshtein 685 distance from the predicted sequence within the 67-word 686 corpora. 687

Since the articulation is silently produced by the individ-688 ual, there is no ground truth audio to verify the accuracy 689 of the actual articulation (for example, the subject might 690 have started with TUESDAY and then corrected for THURS-691 DAY). Therefore, recognizing the inherent ambiguity and 692 emphasizing the model's ability to closely approximate the 693 intended word, we can classify such predictions as cor-694 rect. We denote word error rate calculated in this man-695 ner as WER*. In this case, the best decoding accuracy on 696 Data SMALL-VOCAB is 91% (WER* of just 9%, as opposed to 697 12% with the previous way of calculating WER). 698

When such predictions are classified as correct, the genera tion accuracy of Data NATO-WORDS significantly increases and
 is summarized in table 8.

B.1. Beam search algorithm

We use an algorithm that performs beam search decoding 705 solely on the CTC output probabilities, without incorporat-706 ing any external language models or prior linguistic knowledge. At each timestep, it evaluates the likelihood of extending existing sequences based purely on the symbol probabili-709 ties provided by the CTC output, maintaining a fixed number 710 of the most probable beams (defined by beam width), and 711 ultimately returns the most likely sequence based on the 712 CTC probabilities. 713

714

Table 8. Best decoding accuracy across all GRU models and model sizes of all four subjects (best value is calculated by averaging over 10 random seeds). When multiple entries in the word corpora share the lowest Levenshtein distance with the predicted sequence and the ground truth is among them, we classify it as a *correct* prediction. Results are for Data NATO-WORDS

Subject	Best accuracy (1 - WER*)
1	54.48%
2	59.45%
3	38.90%
4	56.29%
Average	52.28%

C. Background on Riemannian geometry of SPD matrices

Speech articulation involves the coordinated activation of various muscles, with their activation patterns defined by the functional connectivity of the underlying neuromuscular system. Consequently, EMG signals collected from multiple, spatially separated muscle locations exhibit a time-varying graph structure. Gowda et al. demonstrate that the graph edge matrices corresponding to orofacial movements underlying speech articulation are inherently distinguishable on the manifold of SPD matrices. Through experiments with 16 subjects, they highlight the effectiveness of using SPD manifolds as an embedding space for these edge matrices. Building on this foundation, we investigate the temporal evolution of graph connectivity using edge matrices to enable EMG-to-language translation.

Directly working with SPD matrices using affine-invariant or log-Euclidean metrics (Arsigny et al., 2007) involves computationally expensive operations, such as matrix exponential and matrix logarithm calculations. These operations make mappings between the manifold space and the tangent space, and vice versa, computationally intensive. To address this, Lin proposed methods to operate on SPD matrices using Cholesky decomposition. They established a diffeomorphism between the Riemannian manifold of SPD matrices and Cholesky space, which was later utilized by Jeong et al. to develop computationally efficient recurrent neural networks. In Cholesky space, the computational burden is significantly reduced: logarithmic and exponential computations are restricted to the diagonal elements of the matrix, making them element-wise operations. Additionally, the Fréchet mean is derived in a closed form.

Given a set of SPD edge matrices $\mathcal{E}(\tau)$ over different time windows τ , we first calculate their corresponding Cholesky decompositions $\mathcal{L}(\tau) = \text{CHOLESKY}(\mathcal{E}(\tau))$, such 715 that $\mathcal{E}(\tau) = \mathcal{L}(\tau)\mathcal{L}(\tau)^T$. Then, the Fréchet mean of the 716 Cholesky decomposed matrices $\mathcal{L}(\tau)$ is given by

$$\begin{aligned} \mathcal{F}_{\text{CHOLESKY}} &= \frac{1}{n} \sum_{i=1}^{n} \lfloor \mathcal{L}(\tau_i) \rfloor &+ \\ & \exp\left(\frac{1}{n} \sum_{i=1}^{n} \log(\mathbb{D}(\mathcal{L}(\tau_i)))\right) \end{aligned}$$

The Fréchet mean \mathcal{F} on the manifold of SPD matrices is calculated as

$$\mathcal{F} = \mathcal{F}_{\text{CHOLESKY}} \mathcal{F}_{\text{CHOLESKY}}^T$$

In the above equation, $\lfloor \mathcal{L}(\tau) \rfloor$ is the strictly lower triangular part of the matrix $\mathcal{L}(\tau)$, and $\mathbb{D}(\mathcal{L}(\tau))$ is the diagonal part of the matrix $\mathcal{L}(\tau)$.

GRU_A is a standard GRU (Chung et al., 2014). GRU_B is constructed from GRU_A by replacing the arithmetic operations of GRU_A defined in the Euclidean domain with the corresponding operations on the SPD manifold. Gates of GRU_B as defined by Jeong et al. are given below. Given the sparse SPD edge matrices $\sigma(\tau)$ over different time windows τ , we first calculate their corresponding Cholesky decompositions $l(\tau) = \text{CHOLESKY}(\sigma(\tau))$, such that $\sigma(\tau) = l(\tau)l(\tau)^T$.

Update-gate z_{τ} at time-step τ is

$$z_{\tau} = \text{SIGMOID}(w_{z}\lfloor l_{\tau} \rfloor + u_{z}\lfloor h_{\tau-1} \rfloor + b_{z}) + \\ \text{SIGMOID}(b_{z'}[\exp(w_{z'}\log(\mathbb{D}(l_{\tau})) + u_{z'}\log(\mathbb{D}(h_{\tau-1}))]).$$
(1)

Reset-gate r_{τ} at time-step τ is

$$r_{\tau} = \text{SIGMOID}(w_{r} \lfloor l_{\tau} \rfloor + u_{r} \lfloor h_{\tau-1} \rfloor + b_{r}) + \\ \text{SIGMOID}(b_{r'} [\exp(w_{r'} \log(\mathbb{D}(l_{\tau})) + u_{r'} \log(\mathbb{D}(h_{\tau-1}))]).$$
(2)

Candidate-activation vector \hat{h}_{τ} is

$$\hat{h}_{\tau} = \text{TANH}(w_h \lfloor l_{\tau} \rfloor + u_h(\lfloor r_{\tau} \rfloor * \lfloor h_{\tau-1} \rfloor) + b_h) + \text{SOFTPLUS}(b_{h'} \exp(w_{h'} \log(\mathbb{D}(l_{\tau})) + u_{h'} \log(\mathbb{D}(r_{\tau}) * \mathbb{D}(h_{\tau-1})))).$$
(3)

Output vector h_{τ} is

$$h_{\tau} = (1 - \lfloor z_{\tau} \rfloor) * \lfloor h_{\tau-1} \rfloor + \lfloor z_{\tau} \rfloor * \lfloor \hat{h_{\tau}} \rfloor + \exp((1 - \mathbb{D}(z_{\tau})) * \log(\mathbb{D}(h_{\tau-1})) + \mathbb{D}(z_{\tau}) * \log(\mathbb{D}(\hat{h}_{\tau}))). \quad (4)$$

In the above equations, $h_{\tau-1}$ is the hidden-state at time-step $\tau - 1$.

In GRU_C, we define an additional implict layer solved using neural ODEs. The dynamics f of EMG data is modeled by a neural network with parameters Θ . The output state h_{τ} is updated as,

$$h_{\tau-1} \leftarrow \text{ODESOLVE}(f_{\Theta}, \widetilde{\text{LOG}}(h_{\tau-1}), (\tau-1, \tau))$$
$$h_{\tau} = \text{GRU}(l_{\tau}, \widetilde{\text{EXP}}(h_{\tau-1})), \quad (5)$$

where \widehat{LOG} is the logarithm mapping from the manifold space of SPD matrices to its tangent space and \widehat{EXP} is its inverse operation as defined by Lin. GRU is a gated recurrent unit whose gates are given by equations 1 - 4.

Previous work by Gowda & Miller demonstrated the effectiveness of SPD matrices in decoding *discrete* hand gestures from EMG signals collected from the upper limb. Furthermore, SPD matrix representations have been extensively utilized to model electroencephalogram (EEG) signals, although they have never been applied to complex tasks such as sequence-to-sequence speech decoding. For example, Barachant et al.; Barachant et al. employed Riemannian geometry frameworks for classification tasks in EEG-based brain-computer interfaces, while Sabbagh et al. developed regression models based on Riemannian geometry for biomarker exploration using EEG data.

The novelty of our work lies in the algebraic interpretation of manifold-valued data through linear transformations, and the development of models for complex sequence-tosequence tasks. This approach moves beyond the conventional applications of classification and regression.

D. $\sigma(\tau)$ are sparse matrices

770

788

789

790

792

796

797

798

799

800

801

802

803

804

805

806

821

822

823

824

We show that $\sigma(\tau)$ are indeed sparse matrices in figure 7.



Figure 7. Blue: Average value of $\frac{\max(ABS((NON DIAG(\Sigma(\tau))))}{\max(DIAG(\Sigma(\tau)))}$ for all τ in train-validation-test set. Red: Average value of $\frac{\max(ABS((NON DIAG(\sigma(\tau)))))}{\max(DIAG(\sigma(\tau)))}$ for all τ in train-validation-test set. As we can see, $\sigma(\tau)$ are approximately diagonal compared to $\Sigma(\tau)$. We use sparse matrices $\sigma(\tau)$ for EMG-to-language translation. Subjects are abbreviated with notation S1, S2, S3, S4.

E. Effect of training data size on decoding accuracy

We train the GRU_A model with varying train dataset sizes for Data $_{\text{SMALL-VOCAB}}$ and present the decoding accuracy and Levenshtein distance in figure 8. As we can see, decoding accuracy demonstrates a plateauing trend with increase in train dataset size, but importantly, has not saturated yet. In future, we would like to explore if more training data can lead to better decoding accuracy.



Figure 8. Decoding accuracy and Levenshtein distance versus training dataset size for Data $_{\text{SMALL-VOCAB}}$. All experiment parameters are same as in table 1 except for the varying train set size. We use GRU_A for training.

F. Electrode position versus decoding accuracy

The form factor of an EMG-based neuroprosthesis plays a critical role in its usability, particularly in facilitating ease of application and removal. Here, we evaluate three electrode configurations. We have 31 electrodes placed on the throat, neck and the left cheek. The configurations are defined as follows:

① Configuration_A: We consider electrodes placed on the throat and left neck only (10 electrodes).

(2) Configuration_B: We consider electrodes placed on the left cheek, with the neck electrodes excluded (11 electrodes).

(4) Configuration_C: We consider electrodes on the throat and left neck and cheek, excluding those on the right neck (22 electrodes).

This exploration aims to assess the practicality and performance of each configuration to inform design choices for an optimal neuroprosthetic interface. Electrode placement on the throat, the left neck and left cheek is same as described in Gowda et al. Electrode placement on the right neck is symmetrical to that of left neck. Decoding accuracy for various configurations are shown in table 9. The training paradigm is same as in table 1, except for the varying number of electrodes. Decoding accuracy are obtained using GRU_A for Data _{SMALL-VOCAB}.

Table 9. Decoding accuracy with various electrode configurations.

Configuration	Accuracy (1 - WER)
Configuration _A	86.00%
Configuration _B	84.24%
Configuration _C	85.96%

Above results show that EMG based neuroprosthesis can have a small form factor (such as neck only or cheek only), and still provide good decoding accuracy (decoding accuracy using all 31 electrodes is 88%).

G. Text to personalized audio synthesis

The generated personalized audio files will be made available as part of the open-sourced codes.

We synthesize constructed phoneme sequences into personalized audio using methods described by Choi et al. For this, we train the model proposed by Choi et al. on speech corpora provided by Panayotov et al. (LibriSpeech TRAIN-CLEAN-360 and TRAIN-CLEAN-100) and Veaux et al.

8	825 826 827	(VCTK corpus). For few-shot learning, we use a 40-second reference audio clip from the subject (Data $_{\text{SMALL-VOCAB}}$) to capture the speaker's vocal characteristics.
5	328	
8	329	The process involves converting the predicted text into audio
8	330	using Google Text to Speech (aTTS). The aTTS generated
8	331	using Google Text-to-Speech (g115). The g115-generated
8	332	audio is then personalized using the model by Choi et al.,
8	333	leveraging the 40-second reference audio data (the refer-
8	334	ence audio includes inguistic content L that is absent from
8	335	the Data _{SMALL-VOCAB}). This approach ensures that the syn-
8	336	thesized audio closely mimics the speaker's unique vocal
8	337	auridules.
8	338	
8	339	
8	340	
8	341	
8	342	
8	343	
8	344	
8	345	
8	346	
8	347	
8	348	
8	349	
8	350	
8	351	
8	352	
8	353	
8	354	
8	355	
8	356	
8	357	
8	358	
8	359	
8	360	
8	361	
8	362	
8	363	
8	364	
8	365	
8	366	
8	367	
8	368	
8	369	
8	370	
8	371	
8	372	
8	373	
8	374	
8	375	
8	376	
8	377	
8	378	
8	379	