# A Cross-Modal Approach to Silent Speech with LLM-Enhanced Recognition

Tyler Benster [1]   Guy Wilson [2]   Reshef Elisha [3]   Francis R Willett [2]   Shaul Druckmann [4]

## Abstract

Silent Speech Interfaces (SSIs) offer a noninvasive alternative to brain-computer interfaces for soundless verbal communication. We introduce Multimodal Orofacial Neural Audio (MONA), a system that leverages cross-modal alignment through novel loss functions—cross-contrast (crossCon) and supervised temporal contrast (supTcon)—to train a multimodal model with a shared latent representation. This architecture enables the use of audio-only datasets like LibriSpeech to improve silent speech recognition. Additionally, our introduction of Large Language Model (LLM) Integrated Scoring Adjustment (LISA) significantly improves recognition accuracy. Together, MONA LISA reduces the state-of-the-art word error rate (WER) from 28.8% to 12.2% in the Gaddy (2020) benchmark dataset for silent speech on an open vocabulary. For vocal EMG recordings, our method improves the state-of-the-art from 23.3% to 3.7% WER. In the Brain-to-Text 2024 competition, LISA performs best, improving the top WER from 9.8% to 8.9%. To the best of our knowledge, this work represents the first instance where noninvasive silent speech recognition on an open vocabulary has cleared the threshold of 15% WER, demonstrating that SSIs can be a viable alternative to automatic speech recognition (ASR). Our work not only narrows the performance gap between silent and vocalized speech but also opens new possibilities in human-computer interaction, demonstrating the potential of cross-modal approaches in noisy and data-limited regimes.

[1]Neurosciences PhD Program, Stanford University [2]Department of Neurosurgery, Stanford University [3]Department of Chemical Engineering, Stanford University [4]Department of Neurobiology, Stanford University. Correspondence to: Tyler Benster <tbenst@stanford.edu>, Shaul Druckmann <shauld@stanford.edu>.

## 1. Introduction

Silent Speech Interfaces (SSIs) are a branch of human-computer interaction that offers non-invasive means of non-verbal communication. These interfaces may one day be particularly impactful for individuals with speech impairments and in scenarios where traditional vocal communication is impractical or impossible. Despite progress in SSIs, they face significant challenges in achieving sufficiently high accuracy due to both the absence of phonetic content and limited datasets for training and validation.

This research aims to address these limitations by introducing new methodologies that improve silent speech recognition. We propose Multimodal Orofacial Neural Audio (MONA), a novel approach leveraging cross-modal alignment through two new loss functions—cross-contrastive learning (crossCon) and supervised temporal contrastive learning (supTcon). These functions facilitate the training of a multimodal model capable of harnessing audio-only datasets such as LibriSpeech (Panayotov et al., 2015), previously untapped for silent speech recognition.

Additionally, we incorporate Large Language Model (LLM) Integrated Scoring Adjustment (LISA) to significantly improve recognition accuracy. Our methods collectively aim to reduce the word error rate (WER) in silent speech, which is crucial for the practical applicability of SSIs in real-world scenarios.

By narrowing the performance gap between silent and vocalized speech, MONA LISA may help create viable SSI alternatives to existing automatic speech recognition systems. This advancement could enable communication methods for individuals with speech disorders and create a new interface for conversational AI powered by silent speech.

## 2. Related work

SSIs have historically faced challenges such as the absence of phonetic information generated by speech articulators in unrecordable locations and the paucity of training data, impeding their ability to achieve error rates suitable for practical use. Early efforts in SSIs date back to the 1980s, demonstrating successful decoding across a broad array of phonemes with a limited vocabulary using finite automaton (Sugie & Tsunoda, 1985) or maximum likelihood estima-

tion (Morse & O'Brien, 1986). Progress in the field saw a significant leap in the early 2000s with multisession decoding that achieved 87% accuracy on a vocabulary of 10 words using HMMs (Maier-Hein et al., 2005). Jou et al. bootstrapped a silent speech HMM using ASR to 70% accuracy on a 100-word vocabulary (Jou et al., 2006; Jou & Schultz, 2008). Elmahdy & Morsy (2017) used a deep learning model with convolution, RNN, and CTC loss to achieve 20% WER on a 20-word vocabulary. Recently, Gaddy (2022) achieved a breakthrough in open-vocabulary decoding, training a ResNet-Transformer model on an 18-hour dataset to achieve a 29% WER when predicting text in an open-vocabulary setting, and a 36% WER when directly synthesizing audio (Gaddy & Klein, 2020; 2021). Ren et al. (2023) uses the same dataset and architecture to develop an active-learning paradigm to reduce the labeling burden on EMG data collection.

Technological approaches to SSIs have been diverse, including brain implants (Willett et al., 2023; Metzger et al., 2023), lip reading (Shi et al., 2022), ultrasound (Kimura et al., 2019; Sun et al., 2024), MRI (Tang et al., 2023), fNIRS (Liu & Ayaz, 2018), MEG (Défossez et al., 2023), EEG (Lopez-Bernal et al., 2022), radar (Wagner et al., 2022), strain sensors (Kim et al., 2022) and non-audible murmur (Nakajima et al., 2006). Among noninvasive techniques, lip reading and surface electromyography (EMG) are notable for their ability to perform high-accuracy open-vocabulary decoding. Lip reading currently shows the best performance in open vocabulary settings when trained extensively, although its accuracy decreases with reduced training data, underperforming EMG trained on fewer hours of data (Shi et al., 2022; Gaddy, 2022). In theory, EMG may have a lower error floor, as non-visible information can be recorded. When placed over facial muscles, EMG can reliably detect activity related to speech articulation (Schultz & Wand, 2010), and when placed on the throat, EMG can detect internal motion of the larynx and vocal cord (Bruder & Wöllner, 2019).

Substantially more effort has been devoted to the development of machine learning in automatic speech recognition (ASR) than SSIs, so we look to the machine learning ASR literature for inspiration. Initial progress in ASR came from advances in algorithms, from the introduction of beam search (Lowerre, 1976) and hidden Markov models (HMM) (Baker, 1973), to neural networks (Hinton et al., 2012) and end-to-end deep learning (Hannun et al., 2014). Recently, advances in ASR have predominantly come from new loss functions like InfoNCE (van den Oord et al., 2018), unsupervised pretraining as in wav2vec 2.0 (Baevski et al., 2020), and leveraging massively more training data as in Whisper (Radford et al., 2022). By combining contrastive loss functions for unsupervised training and supervised training on LibriSpeech, w2v-BERT (Chung et al., 2021) achieved a record 2.3% WER on within-dataset testing.

## 3. Problem statement

SSIs offer novel communication abilities for people with speech impairments and users in environments where vocal communication is not feasible. SSIs have the potential to restore natural speech in patients with laryngectomy (Sugie & Tsunoda, 1985) or dysarthria and to facilitate seamless and private communication with AI assistants (Kapur et al., 2018). However, these interfaces face inherent challenges due to lack of intelligible sound production. This requires advanced machine learning systems capable of solving a recognition problem that exceeds human capabilities.

The performance threshold for SSIs to become a viable alternative to existing automatic speech recognition (ASR) systems is approximately 15% WER (Pandey et al., 2021). Despite advances in the field, the challenge remains to improve the accuracy and applicability of SSIs to reach this critical performance threshold. Achieving this level of accuracy is crucial for the advancement of SSI technology and unlocking its potential in a wide range of applications, including silent communication in environments sensitive to noise. Our research focuses on EMG data, given its potential for lower error rates and its ability to record non-visible information related to speech articulation (Schultz & Wand, 2010; Bruder & Wöllner, 2019).

## 4. Approach

We introduce two new loss functions for cross-modal contrastive learning, a new latent space alignment approach for silent and vocalized speech using dynamic time warping, and a novel post-processing step following beam search to synthesize the predicted sentence from the top k candidates.

To evaluate our contributions empirically, we utilize the dataset from Gaddy & Klein (2020). The core challenge addressed by Gaddy and Klein was to convert silently mouthed words into audible speech based on EMG sensor data, and so the dataset comprises comprises EMG sensor measurements captured on the face and neck during both vocalized and silently articulated speech. Our methodology closely aligns with the architecture and train / val / test split of Gaddy (2022), ensuring consistency in data processing and model evaluation, allowing us to build upon and extend their work in silent speech recognition.

To augment our dataset size, we use LibriSpeech clean + other for training and LibriSpeech clean for validation and test. We remove all chapters from *The War of the Worlds* by H.G. Wells and *The Adventures of Sherlock Holmes* by Arthur Conan Doyle to avoid test label leakage, as these two books are used in our EMG dataset.

Although we focused our efforts on audio and EMG data for this paper, MONA is readily applicable to any pair or

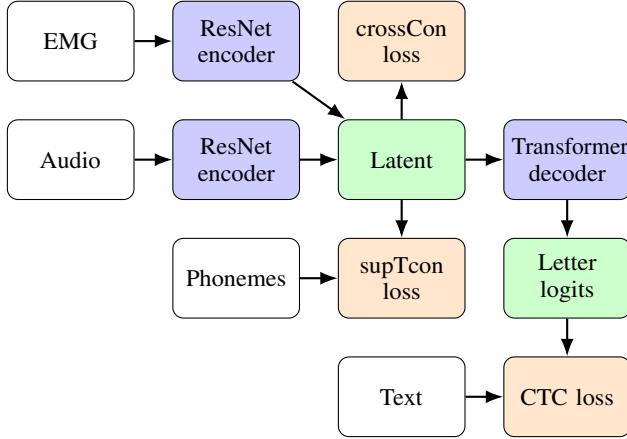group of speech modalities, and LISA is applicable to any speech-to-text prediction model.



*Figure 1.* MONA architecture.

## 4.1. Cross-contrast (crossCon)

For vocalized utterance $u$ with simultaneous EMG and audio recordings, let the output of the EMG encoder be

$$\mathbf{Z}_{\text{emg},u} = \begin{bmatrix} e_{1,u} & e_{2,u} & \cdots & e_{t,u} \end{bmatrix}$$

with $t$ timesteps of 10ms each. The simultaneous output for the audio encoder is

$$\mathbf{Z}_{\text{audio},u} = \begin{bmatrix} a_{1,u} & a_{2,u} & \cdots & a_{t,u} \end{bmatrix}$$

where $\forall t \forall i (e_{t,u}, a_{t,u})$ are latent embeddings for simultaneously recorded EMG and audio data.

We represent all latent embeddings for the minibatch with $n$ utterances by

$$\mathbf{Z} = \begin{bmatrix} \mathbf{Z}_{\text{emg},1} & \mathbf{Z}_{\text{audio},1} & \cdots & \mathbf{Z}_{\text{emg},n} & \mathbf{Z}_{\text{audio},n} \end{bmatrix}$$

where $\mathbf{Z} \in \mathbb{R}^{F \times L}$ for $F$ features and $L$ total embeddings. To simplify indexing, we define $i$ as the index for a specific emg frame and the function $j(i)$ to return the corresponding audio index for utterance $u$. Then, crossCon is defined by Equation 3 for temperature $\tau$:

$$\text{sim}(\mathbf{z}_i, \mathbf{z}_j; \tau = 0.1) = \exp\left(\frac{\cos(\mathbf{z}_i, \mathbf{z}_j)}{\tau}\right) \quad (1)$$

$$\mathcal{L}_i^{\text{cross}} = -\log\left(\frac{\text{sim}(\mathbf{z}_i, \mathbf{z}_{j(i)})}{\sum_{j \neq i}^{L} \text{sim}(\mathbf{z}_i, \mathbf{z}_j)}\right) \quad (2)$$

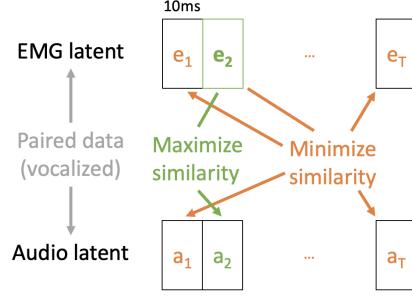$$\mathcal{L}^{\text{cross}} = \frac{1}{L} \sum_{i=1}^{L} \mathcal{L}_i^{\text{cross}} \quad (3)$$



*Figure 2.* **Latent space alignment for time-aligned data with crossCon**. The similarity of latent embeddings for simultaneously recorded EMG and audio data is maximized for the same timestep, and minimized to other timesteps.

Equation 2 can be thought of as a negative log-likelihood loss, where we attempt to classify the specific positive pair of EMG & audio latents among a set of distractors. Equation 3 can be thought of as the average of loss from $L$ classification problems. This differs from CEBRA (Schneider et al., 2023) in that we source distractors from both domains of data. Our formulation encourages the difficult task of creating encoders where cross-modality $e_{t,u}$ and $a_{t,u}$ are more similar than the same-modality $e_{t-1,u}$ and $e_{t,u}$, and therefore encourages a shared latent representation for EMG and audio data (Figure 2).

## 4.2. Supervised temporal contrast (supTcon)

supTcon is suitable when data from different modalities or datasets are not acquired synchronously. We use a class label per latent embedding to align latent representations across data modalities and datasets (Figure 3). Here, the class label is a phoneme or silence as found with Montreal Forced Aligner. To simplify the indexing notation, we introduce the function $p(i)$, which is defined to return the set of all indices corresponding to entries that share the same class label (phoneme) as the entry at index $i$.

$$\mathcal{L}_i^{\text{sup}} = -\sum_{q \in p(i)} \frac{1}{|p(i)|} \log\left(\frac{\text{sim}(\mathbf{z}_i, \mathbf{z}_q)}{\sum_{j \neq i}^{L} \text{sim}(\mathbf{z}_i, \mathbf{z}_j)}\right) \quad (4)$$

$$\mathcal{L}^{\text{sup}} = \frac{1}{L} \sum_{i=1}^{L} \mathcal{L}_i^{\text{sup}} \quad (5)$$

Following the analysis in Khosla et al. (2020), we perform the summation over positives in Equation 4 outside of the log. However, in supTcon, the number of comparisons is quadratic in the total number of timesteps across all examples in the batch, rather than quadratic in the number of examples per batch.

*Table 1.* Overview of datasets, data, and corresponding loss functions. $a$ is an audio utterance, $e$ is an emg utterance, and $y$ is the text label for an utterance. For the Gaddy silent dataset, we have parallel readings of a given utterance ($y_1 = y_2$) under both silent and vocalized conditions. A minibatch may contain examples from multiple datasets, and supTcon is applied to all examples.

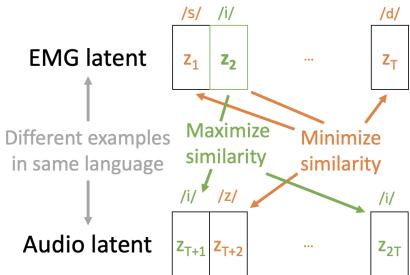| DATASET | DATA | LOSS FUNCTIONS |
|---|---|---|
| GADDY SILENT | $(e_1, y_1), (a_2, e_2, y_2)$ | CROSSCON W/ DTW, SUPTCON W/ DTW, CTC |
| GADDY VOCALIZED | $(a_3, e_3, y_3)$ | CROSSCON, SUPTCON, CTC |
| LIBRISPEECH | $(a_4, y_4)$ | SUPTCON, CTC |



*Figure 3.* **Latent space alignment for independent data with supTcon**. The similarity of latent embeddings with the same phoneme class **/i/** is maximized, and minimized for latent embeddings with different phonemes.

## 4.3. supTcon and crossCon for silent speech with Dynamic Time Warping (DTW)

Since the silent speech data does not contain usable audio, we take advantage of the parallel silent and vocalized utterances in Gaddy & Klein (2020) to align activity across conditions. Since DTW on EMG data performs worse than on audio data (Gaddy, 2022), we instead temporally align the silent condition with the vocalized condition in the latent space. This allows the encoder to act as a denoiser and extract meaningful semantic content that may improve DTW results. Together, we call this approach MONA (Figure 1).

Let the emg latent of silent utterance be $Z_1 \in \mathbb{R}^{F \times T_1}$ with $F$ features and $T_1$ timesteps, and let parallel vocalized emg recording of utterance be $Z_2 \in \mathbb{R}^{F \times T_2}$ and phonemes $P_2 \in \mathbb{P}^{T_2}$ where $\mathbb{P}$ is set that includes phonemes and silence token with $T_2$ timesteps. We first calculate the Euclidean distance matrix $D \in \mathbb{R}^{T_1 \times T_2}$. Then, using DTW, we calculate the alignment of $Z_2$ onto $Z_1$, and use this alignment to create the warped $\hat{Z}_2 \in \mathbb{R}^{F \times T_1}$ and $\hat{P}_2 \in \mathbb{P}^{T_1}$. We can now calculate crossCon and supTcon for the Gaddy silent dataset.

## 4.4. Greedy Class-Weighted Bin Packing

In order to calculate crossCon or supTcon + DTW, we must have at least one silent EMG example with parallel vocalized EMG and Audio per minibatch. As memory usage for contrastive loss functions are quadratic with input length, we must additionally sample batches below a maximum

length to ensure sufficient GPU memory. Finally, since we have massively more audio-to-text data than emg-to-text, we choose to undersample LibriSpeech so that it accounts for approximately 50% of the examples per epoch. Our heuristic for solving this class-proportional bin packing problem is described in Algorithm 1, and is used for our Batch Sampler.

## 4.5. LLM-integrated scoring adjustment (LISA)

During training, we use a 4-gram language model with a beam size of 150 for beam search. By increasing the beam size to 5000 for inference, we reduce WER by 1%. Manual inspection of the top 100 beams revealed that after viewing the list of predictions, in some cases a human could succeed in writing down the correct sentence despite this text not appearing exactly as any specific prediction. As such, rather than calculating the posterior negative log-likelihood (NLL) from the beam search prior and updating with a LLM (Willett et al., 2023), we instead prompt GPT-3.5 (Brown et al. (2020), "gpt-3.5-turbo-16k-0613") or GPT-4 ("gpt-4-0125-preview") with a task description and a list of the top predictions. These predictions can be sourced either from the top $k$ beams from beam search or from an ensemble of models. The latter has a higher diversity of predicted text and is depicted below.

> Your task is to perform automatic speech recognition. Below are multiple candidate transcriptions, listed from most likely to least likely. Choose the transcription that is most accurate, ensuring it is contextually and grammatically correct. Focus on key differences in the options that change the meaning or correctness. Avoid selections with repetitive or nonsensical phrases. In cases of ambiguity, select the option that is most coherent and contextually sound. Respond with the chosen transcription only, without any introductory text.
> after breakfast instead forking at aside to walk down towards the common
> after breakfast stead of working a decided to walk down towards the common
> after breakfast stead working a decided to walk down towards the common

after breakfast instead working a sudden to walk
down towards the common
after breakfast instead of working i decided to
walk down towards the common
after breakfast instead of working a decided to
walk down towards the common
after breakfast instead of working a descended to
walk down towards the common
after breakfast instead of working at a sudden to
walk down towards the common
after breakfast instead of working a decided walk
down towards the common
after breakfast instead of working a decided to
walk down towards the common

This prompt ("Ensemble 10 x top 1") results in the text shown in Table 2, where LISA corrects all four errors present in this utterance.

### 4.6. Experiment setup

We dynamically construct minibatches so that at least one Gaddy silent example (Table 1) is present in each minibatch, and each minibatch is class-balanced for Gaddy silent and Gaddy vocalized, with examples from LibriSpeech filling the remaining 50% of examples. This ensures that each gradient update integrates the loss of each dataset, jointly optimizing for both domains. crossCon and supTcon are quadratic in memory usage in relation to total number of timesteps across the batch, with the latter permitting up to 128K timesteps into the 80GB of VRAM for a Nvidia A100 80GB. With crossCon, up to 256K timesteps were possible on a single GPU.

In an effort to improve numerical stability, we use GeLU activations for the ResNet encoder without pre-norm activation, and, taking inspiration from Balduzzi et al. (2017), we scale the residual path of each block in the ResNet encoder by $\frac{1}{\sqrt{2}}^{\ell}$ where $\ell$ is the layer number. To increase training speed, we leverage the TensorFloat32 format.

We use the same random number generator seed for all data loaders to ensure that each model sees the same sequence of batches during training. An epoch ends when there are no more available samples for a given class, meaning that each epoch typically includes all examples from Gaddy but only a fraction from LibriSpeech. Each model is trained five times with different random initializations. We select the model with the best validation WER on silent EMG for MONA, and ensemble the best 10 models as ranked by validation WER on silent EMG for MONA LISA. For fine-tuning LISA, we split validation into two halves of 100 examples each, fine-tuning on the first 100 and evaluating on the second 100.

A first draft of this manuscript was produced, including all

figures, before the evaluation of the test data. We selected the best MONA and MONA LISA models on the basis of the validation WER for silent EMG. The values reported in Table 3 reflect the test performance of our models, as chosen *a priori* to the first evaluation on the test set.

## 5. Results

To decrease measurement error in comparing model performance, in this section we report the average WER on silent EMG validation from 5 models trained with different initial seeds unless otherwise noted. As shown in Figure 4, our baseline model performs moderately worse on silent speech validation data (30.4% average WER) than the model by Gaddy (2022) (28.8% WER). Beyond our changes to improve numerical stability and training speed, the addition of audio data may additionally negatively impact batch norm statistics in the decoder despite $\lambda_{audio} = 0$ in the loss function (see Appendix A.1), leading to degraded decoder performance. However, our baseline model performs significantly better on the vocal EMG validation data, achieving a 15.1% WER. This improvement may result from the decision to include silent EMG in the training corpus, whereas Gaddy (2022) trains on vocal EMG only to reach 23.3% WER.

We observed higher performance in all model formulations when training on batches with up to 256,000 time steps. However, our implementation of supTcon could only fit on a single A100 with up to 128,000 time steps, so we evaluated most loss functions with this shorter batch length for a fair comparison (Figure 4).

The addition of CTC loss on the Gaddy audio data reduces the silent EMG WER to 27.6%, and incorporating the LibriSpeech Clean & Other datasets result in a significant reduction in WER to 25.5%, surpassing both our baseline model as well as the current state-of-the-art result of Gaddy (2022). The benefits of the EMG & Audio model were more pronounced in the vocal EMG validation set, with the WER dropping to 10.5%.

Latent representations from the Gaddy and LibriSpeech datasets are effectively aligned by crossCon, enhancing the model's robustness in diverse speech scenarios. The integration of crossCon further refined our model's performance, reaching 23.3% WER on silent EMG. Without inclusion of the LibriSpeech dataset, crossCon provides no benefit over the EMG + Audio model (Figure 4). Compared to the baseline model, a model with crossCon trains substantially faster. This also corresponds to improved generalization: at epoch 30, the baseline model has a 58.3% WER on validation data vs 43.4% WER for the model with crossCon (Figure 5). With 256k timesteps per batch, the performance improved further to 22.4%

As crossCon is only applied to the Gaddy dataset, we sought to formulate a contrastive loss function that could directly align EMG and LibriSpeech data. We anticipated that supT-con might therefore improve performance, but found a modest performance degradation in all models trained for this paper (Figure 4).

In order to align silent EMG and audio data, we implemented DTW in latent space such that we can warp parallel vocalized audio latents to silent EMG latents, and apply crossCon and/or supTcon. The combination of crossCon + DTW achieved a 21.4% average WER on silent EMG. We selected the best model, with a silent EMG validation WER of 20.6% as MONA.

The largest improvement came from the addition of more powerful language models. The integration of an LLM to rescore the top 10 beams of the best MONA model ("Direct top 10 beams"; Table 4) corrected for a significant fraction of remaining errors, bringing the validation WER to 18.0%. We noticed that the variability of predictions from beam search was relatively constrained, with most beams differing by only one word, whereas the predictions from models with different initialization varied substantially. By passing the predictions from an ensemble of 10 different MONA models through an LLM, we achieve a 13.1% WER. Finally, by fine-tuning the LLM, we reach the lowest validation WER of 7.3% (Table 4). We call this method of passing an ensemble of 10 MONA models through a fine-tuned LLM, "MONA LISA".

MONA LISA generalizes well to the test set, achieving a record 12.2% WER on silent EMG, and 3.7% WER on vocal EMG (Table 3), reducing the state-of-the-art WER by 57.6% and 84.1% respectively.

Finally, we applied LISA to the Brain-to-Text Benchmark '24. The dataset consists of 12,100 sentences of intended speech by an individual with late-stage Amyotrophic Lateral Sclerosis (ALS), while neural activity was recorded from 256 electrodes in speech-related areas of motor cortex. We fine-tuned LISA on an ensemble of $10 \times$ LSTM, running beam search using a 5-gram language model, using the Pytorch implementation of the model from Willett et al. (2023). At time of writing, LISA is the top-ranked model on the leaderboard, reducing the top WER from 9.8% to 8.9%.

## 6. Discussion

Our approach contrasts with existing methods in SSIs by heavily relying on audio data from vocal utterances to enhance the decoding of silent speech. We also utilize large language models (LLMs) to capitalize on extensive text datasets to improve silent speech recognition. These strategies have significantly narrowed the performance gap between audio and silent sEMG from 26.5% to 9.9% in ab-

Table 2. Example validation transcription before and after LISA

| METHOD | TRANSCRIPTION |
|---|---|
| BEAM SEARCH | AFTER BREAKFAST INSTEAD [OF] **FORKING AT ASIDE** TO WALK DOWN TOWARDS THE COMMON |
| LISA | AFTER BREAKFAST INSTEAD OF WORKING I DECIDED TO WALK DOWN TOWARDS THE COMMON |
| ACTUAL | AFTER BREAKFAST INSTEAD OF WORKING I DECIDED TO WALK DOWN TOWARDS THE COMMON |

Table 3. Test-set word error rate (WER) comparison of different models and enhancements on the Gaddy 2020 benchmark.

| MODEL | LIBRI-SPEECH | GADDY AUDIO | VOCAL EMG | SILENT EMG |
|---|---|---|---|---|
| GADDY (2022) | - | - | - | 28.8% |
| GADDY (2022) | - | - | 23.3% | - |
| GADDY (2022) | - | 11.3% | - | - |
| WHISPER V2 | **2.7%** | **2.3%** | - | - |
| MONA | 9.1% | 7.7% | 8.9% | 22.2% |
| MONA LISA | 5.7% | 2.6% | **3.7%** | **12.2%** |

solute percentage points. Notably, our work has achieved a WER below the critical 15% threshold, indicating the viability of SSIs for open-vocabulary applications.

Compared to the baseline model, a model with crossCon trains substantially faster. This also corresponds to an improved generalization: At epoch 30, the baseline model has a 58.3% WER on validation data vs 43.4% WER for the model with crossCon (Figure 5). When developing cross-Con, we also explored other formulations for encouraging alignment. A simple MSE on paired latent embeddings moved latents towards zero while harming performance. If dot product is used instead of cosine for Equation 1, then the magnitude differs for the latents of each domain, and there are no performance gains.

We observed that CTC loss is only loosely predictive of WER; validation CTC loss typically reached its minima around epoch 40-60, while validation WER reached it's minima around epoch 125-200 (Figure 5). This phenomena has been previously observed for Deepspeech2 (Gururani, 2017) and Conformers as implemented in Nvidia NEMO (psydok, 2022), indicating that these model may become more confident in its predictions (and thus have a lower WER), even as the overall probability distribution modeled by the network, reflected in the CTC loss, becomes less aligned with the true distribution of the validation data. This may indicate that overfitting helps the approximate beam search algorithm find a likely answer with higher reliability. Future research
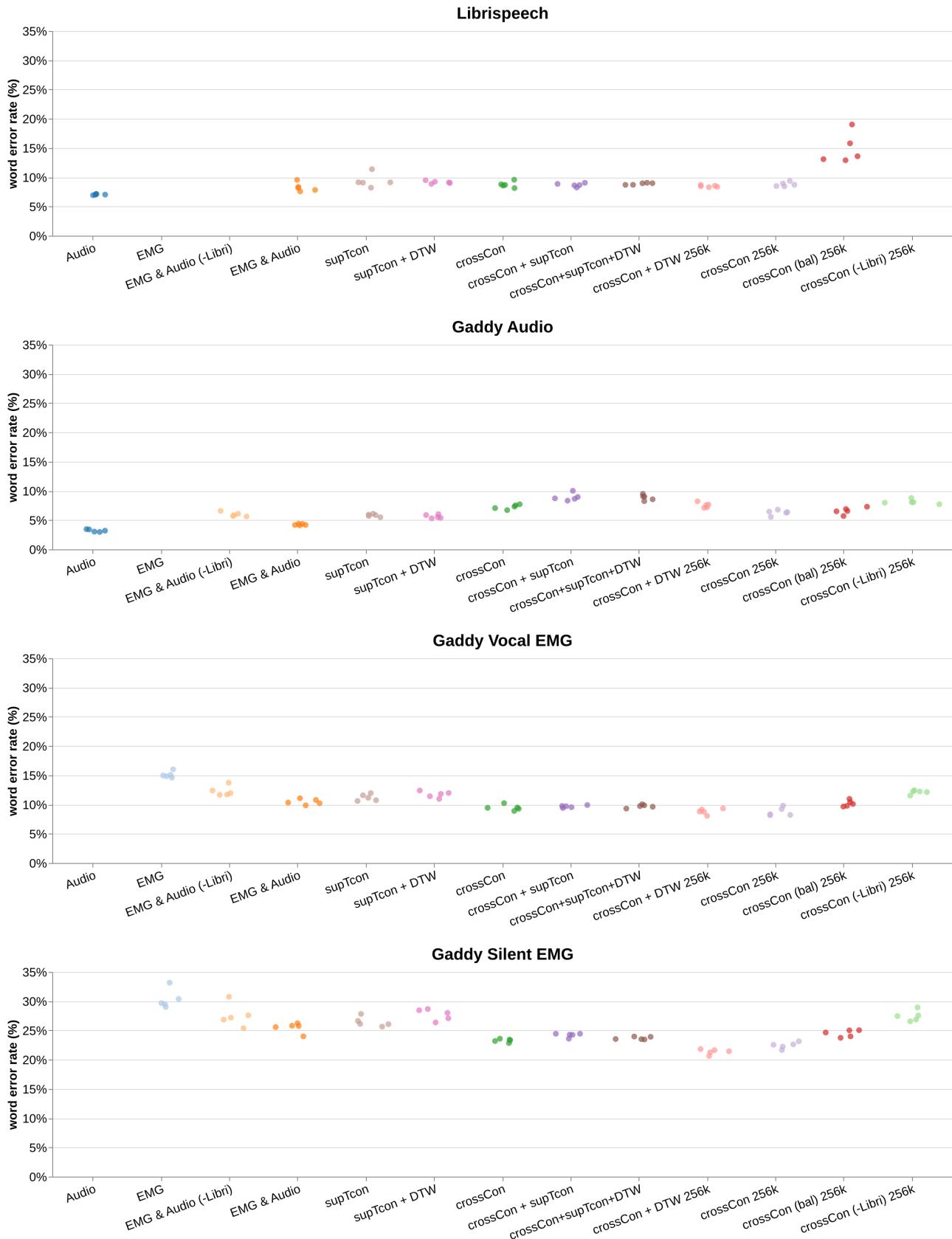
*Figure 4.* Comparison of validation WER for the MONA architecture with varying loss function. -Libri = No LibriSpeech, bal = Balanced audio/EMG sampling.

*Table 4.* Validation WER of LISA approaches on silent EMG data.

| APPROACH | GPT-3.5 | GPT-4 |
|---|---|---|
| CHAIN OF REASONING | 19.5% | 19.2% |
| DIRECT TOP 100 BEAMS | 18.0% | 16.9% |
| DIRECT TOP 10 BEAMS | 18.0% | 17.6% |
| ENSEMBLE 10 × TOP 10 | 13.0% | 14.4% |
| ENSEMBLE 10 × TOP 1 | 13.1% | 13.9% |
| FINE-TUNED 15 × TOP 1 | 7.6% | - |
| FINE-TUNED 10 × TOP 1 | 7.3% | - |

might explore if LISA works better on checkpoints with low CTC loss, or low WER. The overfitting of CTC also invites research that explores new sequence-to-sequence loss functions that generalize without mode collapse.

Although our core contributions held up in the final evaluation on test data, the benefit from fine-tuning on GPT-3.5-turbo generalized poorly on silent EMG data. We tried fine-tuning GPTs on multiple different ensembles: (a) the 10 best models of Figure 4, 5 × (crossCon + DTW 256k) and 5 × (crossCon 256k), (b) we trained an additional 5 × (crossCon + DTW 256k) and selected the 9 × (crossCon + DTW 256k) and 1 × (crossCon 256k) with the lowest validation WER, and (c) the top 15 models 10 × (crossCon + DTW 256k) and (5 × crossCon 256k). Since the fine-tuned model on (a) had the lowest WER on the witheld validation utterances, we selected that model for final test set evaluation per our selection criteria.

In Appendix A.3, we examine the performance of these ensembles as well as other fine-tuning methods, and speculate that alternate construction and selection criteria may be warranted when fine-tuning an LLM for LISA. Fine-tuned ensemble (b) achieved an 8.0% WER on silent EMG test data, and fine-tuned ensemble (c) achieved a 9.1% WER, while fine-tuned ensemble (a) recorded a 12.2% WER. Fine-tuning an ensemble on audio predictions from Librispeech also resulted in <10% WER. We hypothesize that fine-tuning on a diverse dataset may encourage task performance through ensemble weighting, whereas fine-tuning within-domain on limited examples carries the risk of overfitting to the lexicon of the provided predictions.

These learnings informed our training approach for LISA on the Brain-to-text '24 competition. We fine-tuned the 10 × top 1 model on all 600 examples of the "test" data (we consider this the validation set) before evaluating on the 1200 utterances in the "competition" data (we consider this the test set). LISA achieves 13.8% WER on validation prior to fine-tuning, and 10.4% WER on validation after fine-tuning (train / evaluation on same samples). The fine-tuned model generalizes well, reaching a record 8.9% WER on the held-out competition data.

One challenge for deploying LISA in real-time inference settings is the 10x increase in compute required to obtain ensemble predictions. Future work might explore ensemble approximation methods like sampling multiple predictions from a single model with dropout (Gal & Ghahramani, 2016), or acquiring multiple predictions using a mixture of experts (Jacobs et al., 1991). One additional challenge is the instability of ChatGPT API results overtime, which may require new prompt engineering to maintain performance. Future work might explore the use of an open source model such as LLaMA 2 (Touvron et al., 2023) or Mixtral 8x7B (Jiang et al., 2024) for stable performance with long-term reproducibility.

During development, supTcon provided modest improvement under a different set of hyperparameters than ulitimately used in this paper. The approach may warrant continued exploration due to its potential to be used when only silent EMG data are available. Here, we only consider supervised learning where text labels are available; however, supTcon might be extended to the semi-supervised learning case where only subsets of class labels are available, perhaps by using gumbel softmax for phoneme classification instead of DTW. This allows for training on additional data, which may be particularly useful when SSIs are worn outside of the lab.

Inclusion of a text modality during training may further boost performance, but requires a phoneme duration prediction or other alignment approach, similar to the duration / pitch predictor in NaturalSpeech 2 (Shen et al., 2023). Additional context in the form of longer utterances (e.g. 30s) or text transcripts of previous utterances may further allow for LISA to improve imputation of missing phonetic content.

## 7. Conclusion

The present study demonstrates the effectiveness of cross-modal training through novel loss functions and a new latent space alignment approach. By further leveraging LLMs for scoring adjustment, we have significantly reduced the WER in silent speech recognition. Future work might explore the application of these techniques to additional speech modalities, such as invasive BCI or next-generation SSIs.

Our work significantly narrows the performance gap between silent and vocalized speech. This not only illustrates the feasibility of high-accuracy SSIs but also opens new avenues in human-computer interaction, particularly for individuals with speech impairments and in scenarios where traditional vocal communication is impractical. The findings suggest that the performance gap between silent EMG data and ASR for open vocabulary may yet be closed for a single speaker with sufficient electrode coverage or a combination of multiple SSI modalities.

**Algorithm 1** Greedy Class-Weighted Bin Packing
___
   **Input:** item lengths $L$, item class labels $C$, max bin length $M$, class proportions $P$, required classes per bin $I$
   **Output:** list of bins with item indices $B$

   **Initialize** $B$ and sum list $sums$ to empty
   **Group** indices by class into $idx$
   **Shuffle** $idx$ per class
   **Initialize** $debt$ for each class to 0
   **while** num($idx$) $> 0$ for all classes **do**
      **Sample** class $c$ based on $P$
      **if** $debt$ of $c > 0$ **then**
         **Decrement** $debt$ for $c$
         **Continue**
      **end if**
      **Pop** item with length $\ell$ from $L$ and class $c$ from $C$
      **if** there exists a bin where $\ell + \text{sum(lengths in bin)} \leq M$ **then**
         **Find** the first such bin
         **Add** item $idx$ to this bin and update $sums$
      **else**
         **Create** new bin; **Add** item $idx$ to bin
         **if** $I$ classes are missing **then**
            **Add** $idx$ for missing $I$ classes to new bin
            **Increment** $debt$ for added classes
         **end if**
      **end if**
   **end while**
   **Initialize** failure count $f$ to 0
   **Set** failure threshold $T$
   **while** failure count $< T$ and remaining items exist **do**
      **Adjust** $P$ to avoid sampling classes without items
      **Sample** class $c$ based on $P$
      **Try** to add item to existing bin without exceeding $M$
      **if** added successfully **then**
         **Reset** $f$
      **else**
         **Increment** $f$
      **end if**
   **end while**
   **Discard** any remaining items
___

## Software and Data

All software is available at https://github.com/tbenst/silent_speech. Data is available at https://zenodo.org/records/4064409. Training logs and metrics are available for all runs at: https://app.neptune.ai/o/neuro/org/Gaddy

## Impact Statement

This paper presents work with the aim of advancing machine learning techniques for SSIs. Potential societal conse-

quences of our work are numerous, including voice restoration for patient populations with speech disorders, invisible computer interaction, and the decoding of subvocalizations. Decoding inner speech does not appear to be possible with EMG (Nalborczyk et al., 2020), so the prospects of using this work to coercively record private thoughts appear remote. Therefore, similar ethical considerations apply as with other speech decoding technology, like ASR.

## Conflict of Interest

## Acknowledgements

## References

Baevski, A., Zhou, H., Mohamed, A., and Auli, M. Wav2vec 2.0: A framework for self-supervised learning of speech representations. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, NIPS'20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.

Baker, J. K. Machine-aided labeling of connected speech. In *Working Papers in Speech Recognition XI, Technical Reports*, Pittsburgh, PA, 1973. Computer Science Department, Carnegie-Mellon University.

Balduzzi, D., Frean, M., Leary, L., Lewis, J. P., Ma, K. W.-D., and McWilliams, B. The shattered gradients problem: If resnets are the answer, then what is the question? In Precup, D. and Teh, Y. W. (eds.), *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learn-*

*ing Research*, pp. 342–350. PMLR, 06–11 Aug 2017. URL https://proceedings.mlr.press/v70/balduzzi17b.html.

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., Agarwal, S., Herbert-Voss, A., Krueger, G., Henighan, T., Child, R., Ramesh, A., Ziegler, D., Wu, J., Winter, C., Hesse, C., Chen, M., Sigler, E., Litwin, M., Gray, S., Chess, B., Clark, J., Berner, C., McCandlish, S., Radford, A., Sutskever, I., and Amodei, D. Language models are few-shot learners. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems*, volume 33, pp. 1877–1901. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper_files/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf.

Bruder, C. and Wöllner, C. Subvocalization in singers: Laryngoscopy and surface emg effects when imagining and listening to song and text. *Psychology of Music*, 49 (3):567–580, November 2019. ISSN 1741-3087. doi: 10.1177/0305735619883681. URL http://dx.doi.org/10.1177/0305735619883681.

Chung, Y.-A., Zhang, Y., Han, W., Chiu, C.-C., Qin, J., Pang, R., and Wu, Y. W2v-bert: Combining contrastive learning and masked language modeling for self-supervised speech pre-training, 2021.

Défossez, A., Caucheteux, C., Rapin, J., Kabeli, O., and King, J.-R. Decoding speech perception from non-invasive brain recordings. *Nature Machine Intelligence*, 5(10):1097–1107, October 2023. ISSN 2522-5839. doi: 10.1038/s42256-023-00714-5. URL http://dx.doi.org/10.1038/s42256-023-00714-5.

Elmahdy, M. S. and Morsy, A. A. Subvocal speech recognition via close-talk microphone and surface electromyogram using deep learning. In *2017 Federated Conference on Computer Science and Information Systems (FedCSIS)*, pp. 165–168, 2017. doi: 10.15439/2017F153.

Gaddy, D. *Voicing Silent Speech*. PhD thesis, EECS Department, University of California, Berkeley, May 2022. URL http://www2.eecs.berkeley.edu/Pubs/TechRpts/2022/EECS-2022-68.html.

Gaddy, D. and Klein, D. Digital voicing of silent speech. In Webber, B., Cohn, T., He, Y., and Liu, Y. (eds.), *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 5521–5530, Online, November 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.

445. URL https://aclanthology.org/2020.emnlp-main.445.

Gaddy, D. and Klein, D. An improved model for voicing silent speech. In Zong, C., Xia, F., Li, W., and Navigli, R. (eds.), *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, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-short.23. URL https://aclanthology.org/2021.acl-short.23.

Gal, Y. and Ghahramani, Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning, 2016.

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 - ICML '06*, ICML '06. ACM Press, 2006. doi: 10.1145/1143844.1143891. URL http://dx.doi.org/10.1145/1143844.1143891.

Gururani, S. Validation loss increasing while wer decreases. https://github.com/SeanNaren/deepspeech.pytorch/issues/78, 2017. deepspeech.pytorch GitHub issue #78.

Hannun, A. Y., Case, C., Casper, J., Catanzaro, B., Diamos, G., Elsen, E., Prenger, R., Satheesh, S., Sengupta, S., Coates, A., and Ng, A. Y. Deep speech: Scaling up end-to-end speech recognition. *CoRR*, abs/1412.5567, 2014. URL http://arxiv.org/abs/1412.5567.

Hinton, G., Deng, L., Yu, D., Dahl, G. E., Mohamed, A.-r., Jaitly, N., Senior, A., Vanhoucke, V., Nguyen, P., Sainath, T. N., and Kingsbury, B. Deep neural networks for acoustic modeling in speech recognition: The shared views of four research groups. *IEEE Signal Processing Magazine*, 29(6):82–97, 2012. doi: 10.1109/MSP.2012.2205597.

Jacobs, R. A., Jordan, M. I., Nowlan, S. J., and Hinton, G. E. Adaptive mixtures of local experts. *Neural Computation*, 3(1):79–87, 1991. doi: 10.1162/neco.1991.3.1.79.

Jiang, A. Q., Sablayrolles, A., Roux, A., Mensch, A., Savary, B., Bamford, C., Chaplot, D. S., de las Casas, D., Hanna, E. B., Bressand, F., Lengyel, G., Bour, G., Lample, G., Lavaud, L. R., Saulnier, L., Lachaux, M.-A., Stock, P., Subramanian, S., Yang, S., Antoniak, S., Scao, T. L., Gervet, T., Lavril, T., Wang, T., Lacroix, T., and Sayed, W. E. Mixtral of experts, 2024.

Jou, S.-C. S. and Schultz, T. Ears: Electromyographical automatic recognition of speech. In *International Conference on Bio-inspired Systems and Signal Processing*,

2008. URL https://api.semanticscholar.org/CorpusID:5092817.

Jou, S.-C. S., Schultz, T., Walliczek, M., Kraft, F., and Waibel, A. H. Towards continuous speech recognition using surface electromyography. In *Interspeech*, 2006. URL https://api.semanticscholar.org/CorpusID:389078.

Kapur, A., Kapur, S., and Maes, P. Alterego: A personalized wearable silent speech interface. In *23rd International Conference on Intelligent User Interfaces*, IUI'18. ACM, March 2018. doi: 10.1145/3172944.3172977. URL http://dx.doi.org/10.1145/3172944.3172977.

Khosla, P., Teterwak, P., Wang, C., Sarna, A., Tian, Y., Isola, P., Maschinot, A., Liu, C., and Krishnan, D. Supervised contrastive learning. *CoRR*, abs/2004.11362, 2020. URL https://arxiv.org/abs/2004.11362.

Kim, T., Shin, Y., Kang, K., Kim, K., Kim, G., Byeon, Y., Kim, H., Gao, Y., Lee, J. R., Son, G., Kim, T., Jun, Y., Kim, J., Lee, J., Um, S., Kwon, Y., Son, B. G., Cho, M., Sang, M., Shin, J., Kim, K., Suh, J., Choi, H., Hong, S., Cheng, H., Kang, H.-G., Hwang, D., and Yu, K. J. Ultrathin crystalline-silicon-based strain gauges with deep learning algorithms for silent speech interfaces. *Nature Communications*, 13(1), October 2022. ISSN 2041-1723. doi: 10.1038/s41467-022-33457-9. URL http://dx.doi.org/10.1038/s41467-022-33457-9.

Kimura, N., Kono, M., and Rekimoto, J. Sottovoce: An ultrasound imaging-based silent speech interaction using deep neural networks. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*, CHI '19. ACM, May 2019. doi: 10.1145/3290605.3300376. URL http://dx.doi.org/10.1145/3290605.3300376.

Liu, Y. and Ayaz, H. Speech recognition via fnirs based brain signals. *Frontiers in Neuroscience*, 12, October 2018. ISSN 1662-453X. doi: 10.3389/fnins.2018.00695. URL http://dx.doi.org/10.3389/fnins.2018.00695.

Lopez-Bernal, D., Balderas, D., Ponce, P., and Molina, A. A state-of-the-art review of eeg-based imagined speech decoding. *Frontiers in Human Neuroscience*, 16, 2022. ISSN 1662-5161. doi: 10.3389/fnhum.2022.867281. URL https://www.frontiersin.org/articles/10.3389/fnhum.2022.867281.

Lowerre, B. T. The harpy speech recognition system, 1976.

Maier-Hein, L., Metze, F., Schultz, T., and Waibel, A. Session independent non-audible speech recognition using surface electromyography. In *IEEE Workshop on Automatic Speech Recognition and Understanding, 2005.*, pp. 331–336, 2005. doi: 10.1109/ASRU.2005.1566521.

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., Zhuravleva, I., Tu-Chan, A., Ganguly, K., Anumanchipalli, G. K., and Chang, E. F. A high-performance neuroprosthesis for speech decoding and avatar control. *Nature*, 620(7976): 1037–1046, August 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06443-4. URL http://dx.doi.org/10.1038/s41586-023-06443-4.

Morse, M. S. and O'Brien, E. M. Research summary of a scheme to ascertain the availability of speech information in the myoelectric signals of neck and head muscles using surface electrodes. *Computers in Biology and Medicine*, 16(6):399–410, 1986. ISSN 0010-4825. doi: https://doi.org/10.1016/0010-4825(86)90064-8. URL https://www.sciencedirect.com/science/article/pii/0010482586900648.

Nakajima, Y., Kashioka, H., Campbell, N., and Shikano, K. Non-audible murmur (nam) recognition. *IEICE TRANSACTIONS on Information and Systems*, 89(1):1–8, 2006.

Nalborczyk, L., Grandchamp, R., Koster, E. H. W., Perrone-Bertolotti, M., and Lœvenbruck, H. Can we decode phonetic features in inner speech using surface electromyography? *PLOS ONE*, 15(5):e0233282, May 2020. ISSN 1932-6203. doi: 10.1371/journal.pone.0233282. URL http://dx.doi.org/10.1371/journal.pone.0233282.

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)*. IEEE, April 2015. doi: 10.1109/icassp.2015.7178964. URL http://dx.doi.org/10.1109/ICASSP.2015.7178964.

Pandey, L., Hasan, K., and Arif, A. S. Acceptability of speech and silent speech input methods in private and public. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, CHI '21. ACM, May 2021. doi: 10.1145/3411764.3445430. URL http://dx.doi.org/10.1145/3411764.3445430.

psydok. How to interpret the training result: High loss, low wer? https://github.com/NVIDIA/NeMo/discussions/4423, 2022. NVIDIA NeMo GitHub issue #4423.

Radford, A., Kim, J. W., Xu, T., Brockman, G., McLeavey, C., and Sutskever, I. Robust speech recognition via large-scale weak supervision, 2022.

11

Ren, Z., Scheck, K., and Schultz, T. Self-learning and active-learning for electromyography-to-speech conversion. In *15th ITG Conference on Speech Communication*, 10 2023.

Schneider, S., Lee, J. H., and Mathis, M. W. Learnable latent embeddings for joint behavioural and neural analysis. *Nature*, 617(7960):360–368, May 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06031-6. URL http://dx.doi.org/10.1038/s41586-023-06031-6.

Schultz, T. and Wand, M. Modeling coarticulation in emg-based continuous speech recognition. *Speech Communication*, 52(4):341–353, 2010. ISSN 0167-6393. doi: https://doi.org/10.1016/j.specom.2009.12.002. URL https://www.sciencedirect.com/science/article/pii/S0167639309001770. Silent Speech Interfaces.

Shen, K., Ju, Z., Tan, X., Liu, Y., Leng, Y., He, L., Qin, T., Zhao, S., and Bian, J. Naturalspeech 2: Latent diffusion models are natural and zero-shot speech and singing synthesizers, 2023.

Shi, B., Hsu, W.-N., Lakhotia, K., and Mohamed, A. Learning audio-visual speech representation by masked multimodal cluster prediction, 2022.

Sugie, N. and Tsunoda, K. A speech prosthesis employing a speech synthesizer-vowel discrimination from perioral muscle activities and vowel production. *IEEE Transactions on Biomedical Engineering*, BME-32(7):485–490, 1985. doi: 10.1109/TBME.1985.325564.

Sun, X., Xiong, J., Feng, C., Li, H., Wu, Y., Fang, D., and Chen, X. Earssr: Silent speech recognition via earphones. *IEEE Transactions on Mobile Computing*, pp. 1–17, 2024. doi: 10.1109/TMC.2024.3356719.

Tang, J., LeBel, A., Jain, S., and Huth, A. G. Semantic reconstruction of continuous language from noninvasive brain recordings. *Nature Neuroscience*, 26(5):858–866, May 2023. ISSN 1546-1726. doi: 10.1038/s41593-023-01304-9. URL http://dx.doi.org/10.1038/s41593-023-01304-9.

Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., Bikel, D., Blecher, L., Ferrer, C. C., Chen, M., Cucurull, G., Esiobu, D., Fernandes, J., Fu, J., Fu, W., Fuller, B., Gao, C., Goswami, V., Goyal, N., Hartshorn, A., Hosseini, S., Hou, R., Inan, H., Kardas, M., Kerkez, V., Khabsa, M., Kloumann, I., Korenev, A., Koura, P. S., Lachaux, M.-A., Lavril, T., Lee, J., Liskovich, D., Lu, Y., Mao, Y., Martinet, X., Mihaylov, T., Mishra, P., Molybog, I., Nie, Y., Poulton, A., Reizenstein, J., Rungta, R., Saladi, K., Schelten, A., Silva, R., Smith, E. M., Subramanian, R.,

Tan, X. E., Tang, B., Taylor, R., Williams, A., Kuan, J. X., Xu, P., Yan, Z., Zarov, I., Zhang, Y., Fan, A., Kambadur, M., Narang, S., Rodriguez, A., Stojnic, R., Edunov, S., and Scialom, T. Llama 2: Open foundation and fine-tuned chat models, 2023.

van den Oord, A., Li, Y., and Vinyals, O. Representation learning with contrastive predictive coding. *CoRR*, abs/1807.03748, 2018. URL http://arxiv.org/abs/1807.03748.

Wagner, C., Schaffer, P., Amini Digehsara, P., Bärhold, M., Plettemeier, D., and Birkholz, P. Silent speech command word recognition using stepped frequency continuous wave radar. *Scientific Reports*, 12(1), March 2022. ISSN 2045-2322. doi: 10.1038/s41598-022-07842-9. URL http://dx.doi.org/10.1038/s41598-022-07842-9.

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., Shenoy, K. V., and Henderson, J. M. A high-performance speech neuroprosthesis. *Nature*, 620 (7976):1031–1036, August 2023. ISSN 1476-4687. doi: 10.1038/s41586-023-06377-x. URL http://dx.doi.org/10.1038/s41586-023-06377-x.
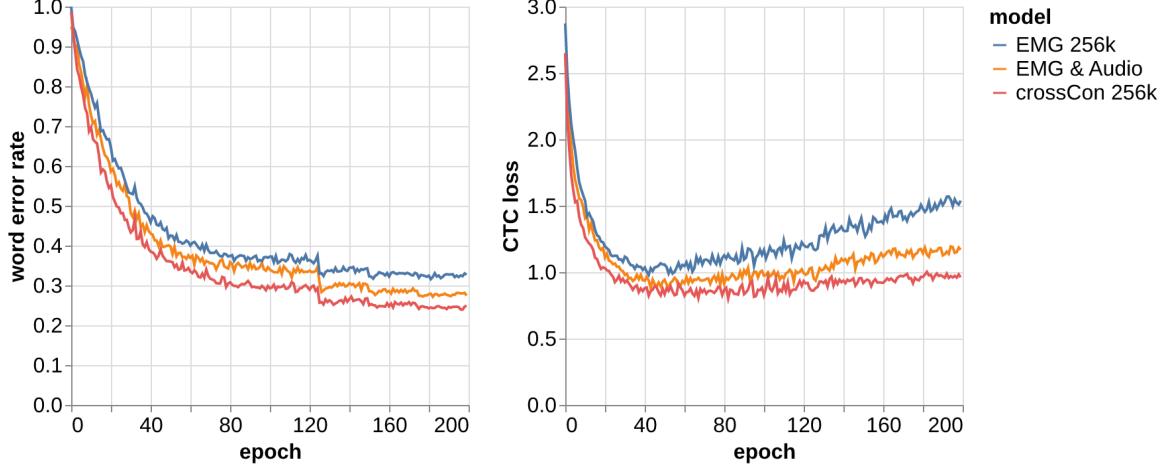
*Figure 5.* Median of five runs WER and CTC loss by training epoch on validation data using 150 beams.

## A. Appendix

### A.1. MONA loss function

Let $\mathcal{L}^{\text{emg}}$ and $\mathcal{L}^{\text{audio}}$ be the Connectionist Temporal Classification (CTC) loss (Graves et al., 2006) for EMG and audio data, respectively. We define the loss for MONA as follows:

$$
\begin{aligned}
\mathcal{L} = \ & \lambda_{\text{emg}} \cdot \mathcal{L}^{\text{emg}} \\
& + \lambda_{\text{audio}} \cdot \mathcal{L}^{\text{audio}} \\
& + \lambda_{\text{cross}} \cdot \mathcal{L}^{\text{cross}} \\
& + \lambda_{\text{sup}} \cdot \mathcal{L}^{\text{sup}}
\end{aligned}
$$

For the experiments in Figure 4, we use the following values of $\lambda$:

*Table 5.* Configuration of lambda parameters for different model runs.

| Model | $\lambda_{\text{audio}}$ | $\lambda_{\text{emg}}$ | $\lambda_{\text{sup}}$ | $\lambda_{\text{cross}}$ |
|---|---|---|---|---|
| Audio | 1 | 0 | 0 | 0 |
| EMG | 0 | 1 | 0 | 0 |
| Audio & EMG | 1 | 1 | 0 | 0 |
| supTcon | 1 | 1 | 0.1 | 0 |
| crossCon | 1 | 1 | 0 | 1 |
| crossCon + supTcon | 1 | 1 | 0.1 | 1 |

### A.2. Balanced EMG/Audio sampling

Both the EMG-only model and the supTcon model suffered from NaN issues during training. We hypothesized that this came from the presence of minibatches with few EMG utterances. To address, we calculated class sampling proportions that in expectation lead to the same number of EMG & Audio utterances per batch with the silent:vocalized EMG dataset ratio: 18.3% silent + parallel EMG, 18.3% LibriSpeech, and 63.3% vocalized EMG (compared with 11.2%, 50%, and 38.8% for all other experiments). This indeed stabilized the gradients for the supTcon model, allowing it to train with (max_len=128000, grad_accum=2) instead of (max_len=128000, grad_accum=4) or (max_len=256000, grad_accum=2), but led to similar validation WER. However, the EMG-only model failed to train with (max_len=128000, grad_accum=2), still requiring the latter two configurations and similarly seeing no change in validation WER. For crossCon, balanced sampling

hurt performance on all tasks. We hypothesize that the training stability benefits of balanced EMG/audio sampling are outweighed by the decrease in performance by undersampling LibriSpeech.

### A.3. Test performance of fine-tuned LISA

*Table 6.* Silent EMG validation and test WER of LISA for the three evaluated ensembles. Ensembles were created and fine-tuned prior to test, and model in bold was chosen based on validation performance.

| | NO FINE-TUNING | | FINE-TUNING | |
| --- | --- | --- | --- | --- |
| ENSEMBLE FOR LISA | VALIDATION | TEST | VALIDATION | TEST |
| **5 × (CROSSCON + DTW), 5 × CROSSCON** | 13.3% | 13.2% | **7.3%** | 12.2% |
| 10 × (CROSSCON + DTW), 5 × CROSSCON | 12.0% | 12.7% | 7.6% | 9.1% |
| 9 × (CROSSCON + DTW), 1 × CROSSCON | 11.9% | 12.8% | 7.6% | 8.0% |

If we fine-tune the first ensemble in Table 6 with ensemble audio predictions from half of the LibriSpeech validation set (270 examples), we achieve a Test performance of 9.2% on silent EMG. Based on these results, we hypothesize that fine-tuning to the silent EMG validation set, where all 100 examples come from *War of the Worlds* may have a higher risk of overfitting to the particular vocabulary in those examples rather than fine-tuning to the larger and more diverse vocabulary in LibriSpeech, which may encourage the LLM to focus on the task of weighting the different models.

Although the fine-tuning results generalize well from dataset-to-dataset for a particular ensemble, they do not generalize well from ensemble-to-ensemble. The test WER is 15.0% using the 10 × ensemble Librispeech fine-tune on the 15 × ensemble predictions (using the same 10 with an additional 5 models), and 21.1% when using a different set of 10 models (row 3 of Table 6, despite the new set consisting of models with lower average validation WER (21.4% vs 21.9%). This supports our hypothesis that fine-tuning on Librispeech is largely learning to weight the predictions of the different models in the ensemble, rather than learning the word statistics of the particular dataset.

### A.4. Alternate LISA prompts

#### A.4.1. CHAIN OF REASONING

> Your task is to perform automatic speech recognition. Below are multiple candidate transcriptions, listed from most likely to least likely. Begin your response with a Chain of Reasoning, explaining your analysis and decision-making process in choosing the most accurate transcription. After your analysis, clearly indicate your final choice with the cue 'TRANSCRIPT: '. Ensure the transcription you choose is contextually and grammatically correct. Focus on key differences in the options that change the meaning or correctness. Avoid selections with repetitive or nonsensical phrases. In cases of ambiguity, select the option that is most coherent and contextually sound. Respond first with your reasoning, followed by 'TRANSCRIPT: ' and then the chosen transcription."

The model exhibited  3% noncompliance to the task on both GPT-3.5 and GPT-4, either refusing to make a selection or failing to respond with "TRANSCRIPT:". We excluded these predictions from Table 4. Without excluding, the WER increases to >30%. The poor performance of chain of reasoning is surprising given the success of this technique in a wide variety of other tasks. We hypothesize that the chain of reasoning may inject the model's own lexicon statistics into the task, and thereby detract from the lexicon statistics of the predictions at hand.

#### A.4.2. NLL LOSS

> Your task is automatic speech recognition. Below are the candidate transcriptions along with their negative log-likelihood from a CTC beam search. Respond with the correct transcription, without any introductory text.

We found that including NLL losses leads to worse performance than a direct approach excluding the phrase "along with their negative log-likelihood from a CTC beam search" and the corresponding NLL values. We hypothesize that the LLM is not capable of doing a Bayesian update by multiplying the token-encoded NLL with its own internal logits and so these numbers are more distracting than simply ranking from best to worst.

## A.5. Test performance of top MONA models

*Table 7.* Validation and Test WER of crossCon+DTW 256k models.

| RUN ID | VALIDATION (%) | TEST (%) |
|---|---|---|
| GAD-984 | 20.63 | 22.17 |
| GAD-992 | 20.79 | 21.69 |
| GAD-986 | 21.26 | 21.75 |
| GAD-996 | 21.32 | 21.87 |
| GAD-993 | 21.45 | 20.72 |
| GAD-987 | 21.45 | 20.90 |
| GAD-988 | 21.63 | 20.96 |
| GAD-995 | 21.69 | 22.54 |
| GAD-983 | 21.82 | 22.11 |
| GAD-994 | 22.27 | 21.87 |
| **AVERAGE** | 21.43 | 21.66 |

To explore to what extent validation performance is predictive of test performance, we conducted a Spearman rank correlation test on Table 7. Spearman's rho was 0.22 (p-val: 0.56) indicating a weekly positive but not statistically significant relationship between validation WER and test WER when choosing between different training runs of the same model architecture differing only by initial seed.