RESOLVING DOMAIN SHIFT FOR REPRESENTATIONS OF SPEECH IN NON-INVASIVE BRAIN RECORDINGS

Anonymous authors

Paper under double-blind review

Abstract

Machine learning techniques have enabled researchers to leverage neuroimaging data to decode speech from brain activity, with some amazing recent successes achieved by applications built using invasive devices. However, research requiring surgical implants has a number of practical limitations. Non-invasive neuroimaging techniques provide an alternative but come with their own set of challenges, the limited scale of individual studies being among them. Without the ability to pool the recordings from different non-invasive studies, data on the order of magnitude needed to leverage deep learning techniques to their full potential remains out of reach. In this work, we focus on non-invasive data collected using magnetoencephalography (MEG). We leverage two different, leading speech decoding models to investigate how an adversarial domain adaptation framework augments their ability to generalize across datasets. We successfully improve the performance of both models when training across multiple datasets. To the best of our knowledge, this study is the first ever application of feature-level, deep learning based harmonization for MEG neuroimaging data. Our analysis additionally offers further evidence of the impact of demographic features on neuroimaging data, demonstrating that participant age strongly affects how machine learning models solve speech decoding tasks using MEG data. Lastly, in the course of this study we produce a new open-source implementation of one of these models to the benefit of the broader scientific community.

031 032 033

034

000

001

002 003 004

006

007

008 009 010

011

013

014

015

016

017

018

019

021

023

024

025

026

027

028

029

1 INTRODUCTION

Applications leveraging recent advancements to decode representations of speech in the brain stand to positively impact the lives of individuals across the world who suffer from impaired verbal communication. While surgically invasive modalities provide the most direct access to the brain, they are practically and ethically prohibitive to conduct at scale. Thus, researchers have increasingly turned towards non-invasive approaches instead. However, non-invasive modalities also come with a unique set of challenges, including a difficult signal-to-noise ratio.

We choose to focus on magnetoencephalography (MEG) as our neuroimaging modality of interest 041 and speech decoding as the principle class of the objectives for the models we train. Specifically, 042 we look at heard speech (listening to someone else speak) decoding as this field is still in its infancy 043 and it is easier to decode than imagined speech (thinking intently of what one is saying without 044 vocalizing it) (Martin et al., 2014). This choice is supported by evidence, albeit contested (Vicente & Langland-Hassan, 2018), of functional overlap between the neural representations of heard and 046 imagined speech (Wandelt et al., 2024). We select MEG because it sits at the intersection between 047 many of the advantages of other non-invasive techniques. While functional magnetic resonance 048 imaging (fMRI) has a stronger spatial resolution than MEG, it is limited in its temporal resolution. Both MEG and electroencephalography (EEG), on the other hand, can record activity on the milisecond timescale at which the brain operates (Hall et al., 2014). Yet in comparison to EEG, MEG has 051 superior spatial resolution, a higher signal-to-noise ratio and a far greater number of scalp-based sensors on average (Hall et al., 2014). In addition, there is evidence to suggest that the use of MEG 052 over EEG is directly correlated with increased performance for speech decoding (Défossez et al., 2023).

054 Overall, the decoding performance of applications using non-invasive modalities continues to lag 055 behind invasive ones - with one reason being the limited scale of the data in most non-invasive stud-056 ies. Despite efforts to acquire increasingly large datasets and curate open neural data repositories, 057 the field has not been able to recreate the the successes that deep learning and "big data" have seen 058 elsewhere. This is due, in large part, to the fact that non-invasive neuroimaging data is inherently difficult to generalize from. For one, different studies employ a myriad of scanners and task designs (Jayalath et al., 2024). Pooling data across scanners and sites then leads to an increase in 060 non-biological variance caused by the differences in the devices and acquisition, including scanner 061 manufacturer (Han et al., 2006)(Takao et al., 2014), upgrade (Han et al., 2006), drift (Takao et al., 062 2011), strength (Han et al., 2006), and gradient nonlinearities (Jovicich et al., 2006). Additionally, 063 within any given study, participants exhibit anatomical and demographic differences that affect the 064 signals recorded from their brains (Jayalath et al., 2024). Providing a reliable means of overcoming 065 this hurdle is an active area of research for both the neuroscience and computer science commu-066 nities. While data harmonization is generally the preferred term among neuroimaging researchers, 067 among computer scientists this problem is most commonly characterized as dataset bias or domain 068 shift (Gretton et al., 2008). In the literature, these two terms are used interchangeably to refer to the 069 same phenomenon. In this study, we present the first application of feature-level harmonization to address domain shift for MEG neuroimaging data. We demonstrate the relevance of this framework for speech decoding by improving the ability of two different networks to generalize across datasets 071 during training to increase performance. 072

073 074

075

2 RELATED WORK

- 076 Domain adaptation (DA) is one approach for solving the domain shift problem which comes from 077 the family of Transfer Learning methods. The bulk of the literature focuses on unsupervised domain adaptation (UDA) as it is a more challenging task that can be trivially adapted for the supervised case. 079 In general, techniques for solving UDA can be categorized as either statistic moment matching (e.g. Long et al. (2018)), domain style transfer (e.g. Sankaranarayanan et al. (2018)), self-training (e.g. 081 Zou et al. (2020); Liu et al. (2021)), or feature-level adversarial learning (e.g. Ganin et al. (2016); He et al. (2020a;b); Liu et al. (2018)) (Liu et al., 2022). Domain shift is often measured by the 083 dissimilarity of the distributions of each domain. A number of metrics have been proposed to this end (Ben-David et al., 2006; 2010; Mansour et al., 2009; 2012; Germain et al., 2013), but the notion most 084 relevant to present study is that of \mathcal{H} -divergence. Based on the work of Kifer et al. (2004), it was later 085 used in the formalization of Ben-David et al. (2006; 2010)'s theory on domain adaptation. This same theoretical framework led to the Domain-Adversarial Neural Networks (DANN) architecture, one of 087 the first successful deep approaches for DA (Ganin et al., 2016). Inspired by generative adversarial 088 networks (GANs), Tzeng et al. (2017) extended the idea behind DANNs with their Adversarial 089 Discriminative Domain Adaptation (ADDA) architecture. An extension of this line of work to Nsource dimensions was subsequently demonstrated by (Zhao et al., 2019). 091
- Given the importance of removing dataset bias (particularly scanner-induced variance) for neu-092 roimaging studies, it is no surprise that a large number of studies have tasked themselves with 093 resolving domain shift in this area. Much of the existing work is based on an empirical Bayes 094 method called ComBat (Johnson et al., 2006). However, ComBat is primarily applied to imagederived values and associations (which MEG is not). The literature in this area has thus focused 096 largely on structural, functional, and diffusion MRI. In fact, DA has been explored more for dif-097 fusion MRI than any other modality - the drawback being many of the methods produced rely on 098 spherical harmonics, limiting the ability to apply them to other neuroimaging techniques (Dinsdale et al., 2021). A few deep approaches have been tried, such as leveraging variational autoencoders (VAEs) (Moyer et al., 2020) and generative models based on the U-Net (Ronneberger et al., 2015) 100 or cycleGAN (e.g. Dewey et al. (2019); Zhao et al. (2019)) architectures. However, these methods 101 are sometimes limited by inherent difficulties validating the harmonized outputs that are generated 102 (Dinsdale et al., 2021). Very few studies to date have leveraged Ben-David et al.'s theoretical UDA 103 framework in the context of neuroimaging data and, to our knowledge, none with MEG data. 104
- 105 An exception to this is the work by Dinsdale et al. (2021). Building on the body of work around 106 \mathcal{H} -divergence, they show that an ADDA-style framework (Tzeng et al., 2017) can be successfully 107 adapted to improve cross-dataset generalization for MRI data. We will refer to this ADDA-style 108 approach as *adversarial harmonization* throughout the remainder of this work. Specific to MEG

108 data, Jayalath et al. (2024) introduce an alternative approach that also manages to find success lever-109 aging data from multiple studies. They propose a pre-training scheme that demonstrates cross-task 110 and cross-dataset generalization wherein the combinations of data used across pre-training and fine-111 tuning encompass different datasets, each employing distinct scanner types and task designs (Jay-112 alath et al., 2024). However, this cross-dataset generalizability remains limited in its efficiency for leveraging aggregated MEG data at the scale required for deep learning. We thus select the architec-113 ture proposed by Jayalath et al. (2024) as one of the two base models we investigate for improving 114 MEG cross-dataset generalization. The second architecture we examine was proposed by Défossez 115 et al. (2023) and reports strong results training over individuals pooled from a single study. Despite 116 showing their model's performance scales with the number of individuals used during training, they 117 do not report a further attempt to train their architecture over multiple datasets. For the sake of 118 convenience, we often refer to these two models by the name of their original code repositories: 119 Brainmagick, for Défossez et al. (2023), and MEGalodon, for Jayalath et al. (2024). 120

121 122

123

124

3 Methods

3.1 DATASETS AND PREPROCESSING

125 This work focuses on four MEG datasets across the two architectures it extends. The Cambridge 126 Centre for Ageing and Neuroscience (Cam-CAN) data repository (Shafto et al., 2014; Taylor et al., 127 2017) contains 641 subjects covering 160 hours of MEG recordings in total. Armeni et al. (2022) 128 contains 30 total hours of recordings (three subjects each listening to 10 hours of speech) and 129 (Gwilliams et al., 2022)'s Manually Annotated Sub-Corpus (MEG-MASC) dataset contains 54 hours 130 (27 subjects each recorded for 2 hours). Lastly, Schoffelen et al. (2019)'s Mother Of Unification Studies (MOUS) consists of 204 subjects recorded for a calculated total of 160 hours. Differences 131 in the devices used during acquisition can carry such a strong signal in the final data that the terms 132 dataset bias and scanner bias are sometimes used interchangeably. However, because the MEG 133 scanner types used are not mutually exclusive among the studies we examine, we make a particular 134 choice to focus on the term dataset bias in a way that is inclusive of, but broader than, acquisi-135 tion device and configuration. We refer to these datasets by the names of their primary authors 136 (i.e. Gwilliams) or their monickers (i.e. MEG-MASC) throughout the rest of this work. Table 1 137 summarizing the above information is included for the reader's convenience. 138

Primary Author	Monicker	Scanner Brand	MEG Hours
Armeni et al. (2022)	-	CTF	30
Gwilliams et al. (2022)	MEG-MASC	KIT	54
Shafto et al. (2014), Taylor et al. (2017)	Cam-CAN	Elekta Neuromag	160
Schoffelen et al. (2019)	MOUS	CTF	160

Table 1: A reference table of the relevant dataset information. Total data volume of each dataset is reported in hours.

147 148

Noting the impact of different demographic distributions between neuroimaging datasets on harmo-149 nization found by Dinsdale et al. (2021), we examine the normalized distributions of both participant 150 age (see Figure 1) and participant sex (see Figure 5 in the Appendix) for each pair of datasets used 151 during training. The Brainmagick experiments leverage the Gwilliams and MOUS datasets, while 152 the MEGalodon experiments use the MOUS and Cam-CAN datasets. We construct a set of subsets 153 both to ease constraints on computational resources as well as examine demographic effects. The 154 Gwilliams and MOUS datasets have relatively equivalent age and sex distributions and therefore no 155 specific measures need to be taken to normalize these features when constructing subsets. While 156 there is also no significant disparity related to the ratios of participant sex between the MOUS and 157 Cam-CAN datasets, we do find a large difference when examining the distributions of participant 158 age. To this end, we construct two different pairs of subsets for these datasets. The first selects indi-159 viduals at random only from the area of overlap between the two age distributions. We refer to these as *balanced subsets*. The second set of subsets, which we call *random subsets*, selects individuals 160 randomly from each dataset such that they approximate the distribution of participant age from the 161 original studies. An even split of males and females overall is controlled for in both cases. In all

164 visualizations relating to each class of subset can be found in Section A.3 of the Appendix. 165 166 Normalized Age Distributions 167 Mean Age: 33.69 ---168 0.25 MOUS 169 Cam-CAN 170 Gwilliams 171 0.20 172 173 174 Density 0.15 175 176 177 178 0.10 179 181 0.05 183 0.00 185 20 30 40 50 60 70 80 90 186 Subject Age (Years) 187

cases, approximately 15 percent of subjects from each dataset are used in total to ensure the ratio of

total subjects and MEG recording hours between any two datasets is maintained. The demographic

Figure 1: The normalized distributions of the ages of the subjects from the MOUS Schoffelen et al. 188 (2019), Cam-CAN Shafto et al. (2014); Taylor et al. (2017), and Gwilliams et al. Gwilliams et al. 189 (2022) datasets. The density plotted along the y-axis represents the proportion (i.e. relative fre-190 quency) of each category within its respective dataset. The mean age calculated over the three datasets is displayed by the dotted line. 192

193

200

191

162

163

194 We focus on harmonizing the deep representations of speech in the brain at the feature level, and as 195 a result apply only whatever minimal preprocessing of MEG data the original papers when imple-196 menting our chosen architectures. The full scope of the preprocessing steps carried out are detailed 197 in Section A.2 of the Appendix.

199 3.2 DECODING TASKS

For the experiments building off of Défossez et al. (2023)'s Brainmagick model, the decoding task 201 is to predict directly the most probable segment of speech stimulus from the corresponding period of 202 MEG data. For the experiments using Jayalath et al. (2024)'s MEGalodon model, the pre-text objec-203 tives are a set of domain-specific classification tasks proposed in the original paper: band prediction, 204 phase shift prediction, and amplitude scale prediction. The speech decoding tasks selected for the 205 fine-tuning phase are speech detection and voicing classification. The goal for speech detection is 206 to determine whether speech has occurred in the section of continuous auditory stimulus given the 207 corresponding segment of MEG data. This should not be confused with the more trivial task of 208 detecting the onset of speech from rest. In voicing classification, phonemes must be classified as 209 voiced or voiceless from the aligned MEG data.

210 211 212

ADVERSARIAL HARMONIZATION 3.3

213 An architecture augmented with adversarial harmonization is generally composed of a feature extractor (which we will refer to as the encoder block), a label predictor (which we will refer to as the 214 task head), and a domain classifier. All parts of the network are first trained together to convergence 215 in a 'warm up' phase to ensure both that the encoder block is able to produce salient features and that the domain classifier is able to distinguish between them accurately. In the harmonization or
'unlearning' phase, the network is then trained via an iterative training procedure composed of three
steps: (1) optimizing the encoder and task head for the primary task, (2) optimizing the domain classifier to identify the remaining dataset bias, and (3) optimizing the encoder to remove the dataset bias
by confusing the domain classifier (Dinsdale et al., 2021). Each of these steps optimizes a unique
loss function and as (2) and (3) are adversarial in nature, they cannot be updated concurrently. Thus, the result is three iterations for each training batch.

Each of the passes during the adversarial phase of the harmonization framework requires its own 224 optimizer, including the initial warm-up phase, for a total of four optimizers used. We keep the 225 original choice of AdamW (Loshchilov & Hutter, 2019) for control epochs (with all parameters 226 updated by a single optimizer) and at first followed the lead of Dinsdale et al. (2021) in using Adam (Kingma & Ba, 2017) for the remaining optimizers. However, hyperparameter testing revealed that 227 the use of a stochastic gradient descent (SGD) optimizer for the adversarial domain classifier led 228 to smoother training during harmonization for both the domain classifier and encoder. The idea to 229 examine different optimizers was informed by the work of Rangwani et al. (2022) who formalize the 230 idea of smoother convergence through the use of SGD when training the adversarial head. This work 231 further confirms their theory over a new domain. A comparison of the performance during training is 232 shown in Section A.8 of the Appendix. For the work involving the MEGalodon framework we retain 233 the learning rate of 0.000066 from the original study for the warm-up phase and follow Dinsdale 234 et al. (2021) in using a learning rate of 0.00001 during harmonization. Alternative choices for the 235 second learning rate are explored, but no real effect is found. Similarly, for the work related to 236 the Brainmagick base model we retain the original paper's learning rate of 0.0003 for the warm-up 237 phase and before reducing it to 0.00001 for all optimizers during the adversarial stage.

238 239

240

3.4 BRAINMAGICK EXPERIMENTS

For the experiments using Défossez et al. (2023)'s architecture as a base, we begin by training the 241 proposed CLIP model first using the original code released alongside the paper and then again with 242 our own implementation of the base model. Leveraging the existing implementation proved to be 243 an unanticipated challenge, in part due to the use of the Python libraries Flashy and Dora developed 244 internally for Facebook Research and relied upon extensively in their code for the Brainmagick pa-245 per. These tools have limited available documentation and are not widely used by other research 246 groups outside of Facebook Research. One major outcome of the present study is therefore sim-247 ply the creation of an open-source implementation of the Brainmagick architecture relying on the 248 standard Pytorch¹ and Lightning² Python libraries. We then look to continue the work of the orig-249 inal authors and investigate whether their architecture benefits from dataset pooling with minimal 250 alterations. We explore a version of "naive" pooling via a pre-training scheme where we train on 251 each dataset, leveraging the saved weights from a previous training run on the other. For example, 252 we train the model on the MOUS dataset except loading the best (as determined by validation loss) saved weights of the baseline run on the Gwilliams dataset. Next, we train the model from random 253 initialization with the same splits but pooling all of the data. Finally, we look to explore whether 254 augmenting the network with adversarial harmonization offers a boost to performance. In all cases 255 we use a 0.7 : 0.1 : 0.2 train/val/test split. 256

The Brainmagick model is implemented faithfully to the original design, though with steps taken 257 to streamline and simplify the repository as a whole - including a critical bug fix related to sensor 258 labeling. The spatial attention layer is used to transform data from both the Gwilliams and MOUS 259 datasets to a uniform number of output channels (270) allowing for the datasets to be processed 260 together. This layer was reported by the authors as being originally designed to support a cross-261 dataset model, as working with multiple studies requires the ability to generalize over different 262 numbers and locations of sensors. As in Défossez et al. (2023), we use a value of 0.2 for the 263 dropout component. The remainder of the *brain model* remains true to the original design. In the 264 case of adversarial harmonization, we treat the brain model as the *encoder block*. The model used 265 alongside the CLIP loss forms the task head, with the CLIP loss remaining as the task loss for 266 harmonization. Lastly, we add our *domain classifier* at the same level as the task head. Following 267 the iterative training regime established by Tzeng et al. (2015), we use cross entropy loss for the

²⁶⁸ 269

¹https://pytorch.org/

²https://lightning.ai/docs/pytorch/stable/

270 domain classifier and employ the confusion loss first proposed by Tzeng et al. (2015) and leveraged 271 by Dinsdale et al. (2021) in optimizing the encoder block to erase the target bias from our deep 272 feature representations. We offer further details and formal definitions of spatial attention, the CLIP 273 loss, and the confusion loss in Section A.4 of the Appendix. The augmented architecture can be 274 found there in Figure 9 as well. All related code is available in this repository³.

275 276

277

3.5 MEGALODON EXPERIMENTS

278 The backbone of Jayalath et al. (2024)'s architecture is a dataset-conditional layer (which projects all MEG recordings into a shared dimensional space) and cortex encoder (which extracts deep rep-279 resentations of brain activity). Additionally, we choose to follow the original authors in applying 280 the optional subject embeddings to the final output of the backbone. When pre-training, these fea-281 tures have a projection applied to them before being used to solve a series of "pre-text" tasks. If 282 fine-tuning, the output of the encoder is used directly for the speech decoding tasks. In applying 283 adversarial harmonization, we treat all layers up to and including the optional subject conditioning 284 as the encoder block. Similarly, the areas responsible for the pre-text and fine-tuning objectives are 285 grouped as the task head. At the same level as the task head, we then add a domain classifier that 286 forms the adversarial component of our harmonization scheme. We also keep the task loss from 287 the original MEGalodon framework for both phases of harmonization (warm-up and adversarial). 288 As before, we use cross entropy loss and the confusion loss for each backward pass involving the 289 domain classifier. Further details on the original MEGalodon architecture and an illustration of the augmented architecture are given in Section A.5 and Figure 10 of the appendix, respectively. 290

291 The MEGalodon pre-text tasks, by design, apply some transformation to the input and then pass 292 this transformed input through the backbone in order to create a set of transformed features to be 293 used for the actual prediction task. However, because the effect is running the entire encoder and 294 creating a unique set of features one time for each task, we pass each additional feature vector through the domain classifier as well. The losses from each of these are aggregated by summation 295 before performing the backwards pass. As we use three pre-text tasks, every training step handles 296 four total feature vectors. We set the value of α , the scaler variable applied to the cross-entropy loss 297 of the domain classifier, to 0.25 to account for this. 298

299 In the case where participant age is targeted instead of dataset bias, we again follow the harmoniza-300 tion roadmap laid out by Dinsdale et al. (2021). In order to adapt the continuous feature of age to 301 a categorical one such that it can be predicted by the adversarial classifier, we create 72 single-year bins spanning from the youngest age across all the datasets (18) to the oldest (89). However, it is 302 also important to capture the fact that for a true age of 25, a prediction of 24 is more accurate than 303 a prediction of 63. To account for this, both the true age labels and the predicted ages (produced by 304 applying a softmax activation to the output of the classifier and then taking the argmax) are converted 305 to softmax labels normally distributed as a $\mathcal{N}(\mu, \sigma^2)$ where μ is equal to the discrete age value and 306 σ is set to 10. Lastly, the cross-entropy loss is swapped out in favor of the Kullback-Leibler (KL) 307 divergence, where we treat the distance of the softmax distribution of the predicted ages from the 308 softmax distribution of the true ages as the loss value. All code related to the augmented MEGalodon 309 architecture is available in this repository⁴.

- 310 311
- 4 RESULTS
- 312 313

321

322

323

4.1 BRAINMAGICK RESULTS 314

315 The re-implementation of the Brainmagick architecture proposed by Défossez et al. (2023) using 316 more widely documented libraries is a success. Our version performs comparatively to the results 317 reported in the paper, albeit with a slight reduction in performance which we attribute to our choice 318 to train all the runs related to our build on a single GPU. As the original authors note in their 319 repository⁵, the number of GPUs used during training can have a large impact when using contrastive 320 losses and for this reason we use the single GPU results of the control runs of our build as our primary

³https://anonymous.4open.science/r/BMBU-9C3E/README.md

⁴https://anonymous.4open.science/r/megalodon-harmonizer-22A3/README.md

⁵https://github.com/facebookresearch/brainmagick

324	Full-Run Results		Top-10 Accuracy	
325	Mothod Training Dat		Gwilliams MOUS	
326	Methou	II alling Data	Owinianis	10003
327	Control (Official repo)	Gwilliams	70.7%, 70.7%*	-
328	Control (Official repo)	MOUS	-	68.5%, 67.5%*
320	Control (Our implementation)	Gwilliams	69.8%	-
329	Control (Our implementation)	MOUS	-	68.1%
330	Pre-trained on MOUS	Gwilliams	68.8%	-
331	Pre-trained on Gwilliams	MOUS	-	67.1%
332	Control	Gwilliams + MOUS	$68.8\%\pm0.5$	$66.8\%\pm0.4$
333	Harmonized	Gwilliams + MOUS	$\textbf{71.0\%}\pm0.2$	$\textbf{68.6\%} \pm 0.2$

Table 2: Following the convention of Défossez et al. (2023), we report Top-10 segment-level ac-335 curacy with confidence intervals calculated over 3 seeds. Results as reported in the original study 336 are denoted by a single asterisk (*). The best performance recorded over each validation dataset is 337 marked in bold.



360 Figure 2: t-SNE plot (van der Maaten & Hinton, 2008) of the activations of the final layer of the encoder block using the Brainmagick (Défossez et al., 2023) base architecture while training over 361 subsets. The left plot shows the control and the right after harmonization. The control data is 362 approximately linearly separable, while the harmonized data is closely mixed. 363

baseline comparison. We do not find that the base Brainmagick architecture is effective in cross-366 dataset generalization, as performance decreases when using both MOUS and Gwilliams during 367 training. This was true for both the attempt at naively combining the datasets via a pre-training

364

334

368 approach as well as training on the datasets pooled together directly. However, we find that using 369 adversarial harmonization yields a 2.2% increase in performance evaluating over the Gwilliams test 370 split and 1.8% increase in performance for the MOUS test split when pooling datasets for training. 371 We conduct a one-sided independent samples *t*-test using the results collected across three seeds 372 and find that our augmentations are statistically significant (p < 0.05) for both the Gwilliams test 373 split (p = 0.012) and MOUS test split (p = 0.011). In fact, adversarial harmonization allows the 374 architecture to successfully combine the datasets during training to improve top-10 accuracy even 375 over the results reported in the original paper. All results from training with the full datasets are shown in Table 2. These findings are additionally supplemented by analogous results collected over 376 subsets of the data (see Table 4 in the Appendix). Together, these results demonstrate a scaling effect 377 as overall data volume is increased.

378	Armeni Fine-Tuning		Balanced Accuracy	
379 380	Method	Method Pre-training Data		Voicing
381	Control	Balanced (M+CC) Random (M+CC)	57.29% 56.53%	52.60% 52.38%
383	Warm-up Only	Balanced (M+CC)	57.76%	52.35%
384 385	Short Warm-up (dataset) Harmonized (dataset)	Balanced (M+CC) Balanced (M+CC)	56.33% 55.04%	51.99% 52.44%
386	Harmonized (dataset) Harmonized (age)	Random (M+CC) Random (M+CC)	56.25% 50.68%	52.42% 50.82%
387	Harmonized (both)	Random (M+CC)	56.15%	52.65%

389 Table 3: We report the balanced accuracy results for the speech detection and voicing classification 390 tasks, fine-tuning and testing on the Armeni dataset. Balanced refers to subsets with no strong bias related to the age distribution of the participants with respect to either dataset, while for Random 391 this is not controlled for. The confound being harmonized is denoted in parentheses in the method 392 column. All runs reported here are trained for 200 epochs. Short Warm-up denotes pre-training 393 the encoder for 100 epochs without the domain classifier, before training with the classifier for an 394 additional 10 epochs and beginning harmonization at epoch 110. Warm-up only indicates training 395 in the warm-up phase for the entire 200 epochs. 396

397

403 404

We further support our claim of cross-dataset generalization by observing that upon beginning the harmonization phase, the dataset classifier is reduced from an average 99.9% accuracy to an average 79.7% and 67.9% accuracy in the full and subset cases, respectively. Additionally, we extract the features produced by the final layer of the encoder block and visualize the change in the separability of the activations through a t-SNE plot (van der Maaten & Hinton, 2008) shown in Figure 2.

4.2 MEGALODON RESULTS

405 The results support our hypothesis that the skewed demographic features of the MOUS and Cam-406 CAN datasets are a possible cause of the difficulty Jayalath et al. (2024)'s model has attempting 407 to scale the number of datasets used during pre-training. We show that the performance of the 408 original model is improved for both decoding tasks when training with the age-balanced subsets 409 as opposed to the random subsets (see Table 3). We conduct experiments just on the age-balanced 410 subsets to determine the degree to which dataset of origin exists as a confound independent of age 411 distribution. Even in this case, a randomly initialized dataset classifier achieves 99.9% accuracy after 412 only a single epoch of training. The features produced by the pre-training scheme therefore have 413 dataset-identifiable aspects beyond those related to participant age. Augmenting the model with 414 adversarial harmonization, we successfully manage to lower the dataset classification accuracy to 51% on average (a reduction of 48.9%). Using the random subsets and targeting age bias, we see in 415 Figure 4 that the softmax probability distribution of the classifier for participant age is driven closer 416 to universal chance after training with adversarial harmonization. Fine-tuning results for dataset 417 harmonization, age harmonization, and jointly harmonizing for both age and dataset bias are also 418 shown in Table 3. 419

We find that pre-text task validation loss is a direct proxy for speech decoding performance in the case of speech detection, but they become uncoupled for voicing classification. Harmonization has a negative effect on speech detection performance in all cases, yet jointly harmonizing for both age and dataset bias delivers better voicing classification performance than the control on both the random *and* age-balanced subsets - despite having a much worse final pre-text task validation loss at the end of training. t-SNE (van der Maaten & Hinton, 2008) plots comparing the final encoder block activations from the end of the warm-up phase to the end of harmonization are shown in Figure 3.

427

5 DISCUSSION

428 429

The increase in fine-tuning performance for the MEGalodon architecture when using the agebalanced subsets compared to the random subsets (Table 3) demonstrates that the demographic features of a subject strongly affect the characteristics of the data collected using MEG devices.



Figure 4: Softmax values of the true and predicted ages averaged over a single batch and converted into Gaussian distributions. The dashed line represents the value of setting all softmax values equal. The plot on the left shows the output of the model at the end of the warm-up phase, while the one on the right shows the output after an additional 100 epochs of adversarial harmonization. After harmonization, the distribution is flattened towards an equal distribution across all ages.

486 However, the particular direction of this effect is likely due to the fact that the datasets used for eval-487 uation (whether that be Armeni or Gwilliams) both have distributions of participant age much more 488 similar to that of the MOUS dataset. While the Cam-CAN dataset includes individuals ranging from 489 18 years old to 89 years old, the other three datasets don't record any participants older than 41, 490 with most younger than 30. Were the fine-tuning datasets to have more inclusive age-distributions, as is the case with Cam-CAN, the results of the control comparison might be flipped. This indicates 491 a broader need within the neuroimaging community to increase efforts to recruit older participants 492 when conducting these studies. Still, the results found here indicate that harmonization methods 493 could play a critical role in reducing age-related effects in both present and future studies. <u>191</u>

495 Additionally, we find adversarial harmonization can be extremely unstable, with task loss diverging 496 sharply when the harmonization phase begins. This effect can be mitigated for dataset bias through hyperparameter selection. As the models are still able to reduce the task loss after beginning har-497 monization, the dataset-identifying features present in the encoder output must not be necessary to 498 perform speech decoding. We were not able to complete equivalent hyperparameter testing for the 499 experiments targeting age as a confound. While the harmonized Brainmagick model was able to im-500 prove validation loss over the control, the MEGalodon model could not within our training horizon. 501 We believe one of the differentiating factors between these two architectures is the speed at which 502 they reach convergence. This is examined in further detail, including related runtime and batch size 503 experiments, in Section A.9 and A.8 of the Appendix. 504

Lastly, as we note above, augmenting the MEGalodon architecture with adversarial harmonization 505 has a negative effect for speech detection and a positive one for voicing classification. This is likely 506 due to the difference in the fine-tuning protocol established by Jayalath et al. (2024) between the 507 tasks rather than a significant qualitative difference. 'Shallow' fine-tuning the network for speech 508 detection only updates the task head, but both the encoder and task head are updated when 'deep' 509 fine-tuning for voicing classification. The discrepancy in performance between the decoding tasks 510 could indicate that harmonization drives the features into a more universal brain representation but 511 at the initial cost of task-specific (in this case *task* referring broadly to all speech perception) perfor-512 mance. This is resolved in the deep fine-tuning case as the encoder is given the chance to re-focus 513 on the downstream task, but remains salient for shallow fine-tuning as the encoder is kept frozen.

514 515

6 IMPACT AND FUTURE WORK

516 517

533

To the best of our knowledge, this study is the first ever application of feature-level, deep learning based harmonization for MEG neuroimaging data. We demonstrate some of the unique challenges of harmonizing MEG data when compared to other modalities and demonstrate that age-related features strongly affect how machine learning models solve speech decoding tasks from MEG data. Using the Brainmagick and MEGalodon architectures as a base, we achieve success augmenting the ability of two different, leading speech decoding models to generalize between datasets. We produced these results even without extensive hyperparameter testing, meaning there were likely performance gains still left on the table.

525 A continuation of this work involving an exhaustive hyperparameter search for both models and 526 unrestricted training time for the augmented MEGalodon model would be well-warranted, as it 527 could not be completed within the scope of the present study. Future work should also explore how the effects reported here hold when pooling training data from upwards of three datasets. We 528 acknowledge the limitations of this study in the Appendix. However, the scope of the work produced 529 remains significant and we believe this study to be of value to the field. As a whole, the results 530 reported here are evidence for the potential of adversarial harmonization to aid in solving the scaling 531 problem for deep learning applications when it comes to MEG data. 532

- 534 REFERENCES
- Kristijan Armeni, Umut Güçlü, Marcel van Gerven, and Jan-Mathijs Schoffelen. A 10-hour withinparticipant magnetoencephalography narrative dataset to test models of language comprehension. *Scientific Data*, 9, 06 2022. doi: 10.1038/s41597-022-01382-7.
- 539 Alexei Baevski, Henry Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework for self-supervised learning of speech representations, 2020. URL https://arxiv.

540 541	org/abs/2006.11477.
542	Shai Ben-David, John Blitzer, Koby Crammer, and Fernando Pereira. Analysis of representations for
543	domain adaptation. In Proceedings of the 19th International Conference on Neural Information
544	Processing Systems, NIPS'06, pp. 137–144, Cambridge, MA, USA, 2006. MIT Press.
545	
546	Shai Ben-David, John Blitzer, Koby Crammer, Alex Kulesza, Fernando Pereira, and Jennifer Wort-
547	man Vaughan. A theory of learning from different domains. Mach. Learn., $79(1-2):151-175$,
548	may 2010. ISSN 0885-0125. doi: 10.1007/S10994-009-5152-4. UKL https://doi.org/
549	10.1007/510994-009-5152-4.
550	Blake E. Dewey, Can Zhao, Jacob C. Reinhold, Aaron Carass, Kathryn C. Fitzgerald, Elias S.
551	Sotirchos, Shiv Saidha, Jiwon Oh, Dzung L. Pham, Peter A. Calabresi, Peter C.M. van Zijl, and
552	Jerry L. Prince. Deepharmony: A deep learning approach to contrast harmonization across scan-
553	ner changes. Magnetic Resonance Imaging, 64:160–170, 2019. ISSN 0730-725X. doi: https://doi.
554	org/10.1016/j.mri.2019.05.041. URL https://www.sciencedirect.com/science/
555	article/pii/S0730725X18306490. Artificial Intelligence in MRI.
556	Nicola K Dinsdale Mark Jenkinson and Ana II Namburete. Deep learning-based unlearn-
557	ing of dataset bias for mri harmonisation and confound removal NeuraImage 228:117689
558	2021. ISSN 1053-8119. doi: https://doi.org/10.1016/i.neuroimage.2020.117689. URL https://
559	//www.sciencedirect.com/science/article/pii/S1053811920311745.
560	
561	Alexandre Défossez, Charlotte Caucheteux, Jérémy Rapin, Ori Kabeli, and Jean-Rémi King. De-
562	coding speech perception from non-invasive brain recordings. <i>Nature Machine Intelligence</i> , 5
563	(10):1097-1107, October 2023. ISSN 2522-5839. doi: 10.1038/s42256-023-00714-5. URL
564	nttp://dx.dol.org/10.1038/s42256-023-00/14-5.
565	Yaroslav Ganin, Evgeniva Ustinova, Hana Aiakan, Pascal Germain, Hugo Larochelle, François
566	Laviolette, Mario Marchand, and Victor Lempitsky. Domain-adversarial training of neural net-
567	works, 2016. URL https://arxiv.org/abs/1505.07818.
568	
569	Pascal Germain, Amaury Habrard, François Laviolette, and Emilie Morvant. A pac-bayesian ap-
570	proach for domain adaptation with specialization to linear classifiers. In Sanjoy Dasgupta and
571	David McAnester (eds.), Froceedings of the Soin International Conference on Machine Learning, volume 28 of Proceedings of Machine Learning Research pp. 738–746. Atlanta Georgia USA
572	17-19 Jun 2013 PMLR LIRL https://proceedings.mlr.press/v28/germain13
573	html.
574	
575	Anne-Lise Giraud and David Poeppel. Cortical oscillations and speech processing: Emerg-
576	ing computational principles and operations. Nature neuroscience, 15:511-7, 03 2012. doi:
577	10.1038/nn.3063.
578	Alexandre Gramfort Martin Luessi Eric Larson Denis A Engemann Daniel Strohmeier Christian
579	Brodbeck, Roman Goi, Mainak Jas, Teon Brooks, Lauri Parkkonen, and Matti Hämäläinen. Meg
580	and eeg data analysis with mne-python. <i>Frontiers in Neuroscience</i> , 7, 2013. ISSN 1662-453X.
581	doi: 10.3389/fnins.2013.00267. URL https://www.frontiersin.org/journals/
582	neuroscience/articles/10.3389/fnins.2013.00267.
583	
584	Arthur Gretton, Alexander Smola, Jiayuan Huang, Marcel Schmittfull, Karsten Borgwardt, and
585	Definition Scholkopi. Covariate Snift by Kernel Mean Matching, volume 3, pp. 131–160. MIT Press 12 2008 ISBN 9780262170055 doi: 10.7551/mitoress/0780262170055.002.0009
586	11035, 12 2000. ISBN 9700202170055. doi: 10.7551/iiiipic55/9700202170055.005.0008.
587	Laura Gwilliams, Graham Flick, Alec Marantz, Liina Pylkkanen, David Poeppel, and Jean-Remi
588	King. Meg-masc: a high-quality magneto-encephalography dataset for evaluating natural speech
589	processing, 2022. URL https://arxiv.org/abs/2208.11488.
590	Emme I Hall Olda E Dahara Data C Mark and Mark I D 1 701 1 4 1 1
591	Emma L. Hall, Sian E. Robson, Peter G. Morris, and Matthew J. Brookes. The relationship be- twoon mag and fmri. Neurolmaga, 102:80, 01, 2014, ISSN 1052, 8110, doi: https://doi.org/10.
592	1016/i neuroimage 2013 11 005 URL https://www.scioncodiroct.com/scionce/
593	article/pii/S1053811913010975. Multimodal Data Fusion.

594 Xiao Han, Jorge Jovicich, David Salat, Andre van der Kouwe, Brian Quinn, Silvester Czanner, 595 Evelina Busa, Jenni Pacheco, Marilyn Albert, Ronald Killiany, Paul Maguire, Diana Rosas, Nikos 596 Makris, Anders Dale, Bradford Dickerson, and Bruce Fischl. Reliability of mri-derived measure-597 ments of human cerebral cortical thickness: the effects of field strength, scanner upgrade and 598 manufacturer. NeuroImage, 32(1):180-194, August 2006. ISSN 1053-8119. doi: 10.1016/ j.neuroimage.2006.02.051. URL https://doi.org/10.1016/j.neuroimage.2006. 02.051. 600 601 Gewen He, Xiaofeng Liu, Fangfang Fan, and Jane You. Classification-aware semi-supervised do-602 main adaptation. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern 603 Recognition Workshops. IEEE, 06 2020a. doi: 10.1109/CVPRW50498.2020.00490. 604 Gewen He, Xiaofeng Liu, Fangfang Fan, and Jane You. Image2audio: Facilitating semi-supervised 605 audio emotion recognition with facial expression image. In Proceedings of the IEEE/CVF 606 Conference on Computer Vision and Pattern Recognition Workshops. IEEE, 07 2020b. doi: 607 10.1109/CVPRW50498.2020.00464. 608 Dulhan Jayalath, Gilad Landau, Brendan Shillingford, Mark Woolrich, and Oiwi Parker Jones. The 609 brain's bitter lesson: Scaling speech decoding with self-supervised learning, 2024. URL https: 610 //arxiv.org/abs/2406.04328. 611 612 W. Evan Johnson, Cheng Li, and Ariel Rabinovic. Adjusting batch effects in microarray expression 613 data using empirical Bayes methods. Biostatistics, 8(1):118-127, 04 2006. ISSN 1465-4644. 614 doi: 10.1093/biostatistics/kxj037. URL https://doi.org/10.1093/biostatistics/ 615 kxj037. 616 Jorge Jovicich, Silvester Czanner, Douglas Greve, Elizabeth Haley, Andre van der Kouwe, Randy 617 Gollub, David Kennedy, Franz Schmitt, Gregory Brown, James MacFall, Bruce Fischl, and An-618 ders Dale. Reliability in multi-site structural mri studies: Effects of gradient non-linearity cor-619 rection on phantom and human data. NeuroImage, 30(2):436-443, 2006. ISSN 1053-8119. doi: 620 https://doi.org/10.1016/j.neuroimage.2005.09.046. URL https://www.sciencedirect. 621 com/science/article/pii/S1053811905007299. 622 Daniel Kifer, Shai Ben-David, and Johannes Gehrke. Detecting change in data streams. In Proceed-623 ings of the Thirtieth International Conference on Very Large Data Bases - Volume 30, VLDB '04, 624 pp. 180–191. VLDB Endowment, 2004. ISBN 0120884690. 625 Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization, 2017. URL 626 https://arxiv.org/abs/1412.6980. 627 628 Xiaofeng Liu, Yang Zou, Lingsheng Kong, Zhihui Diao, Junliang Yan, Jun Wang, Site Li, Ping 629 Jia, and Jane You. Data augmentation via latent space interpolation for image classification. In 630 24th International Conference on Pattern Recognition (ICPR), pp. 728–733. ICPR, 08 2018. doi: 631 10.1109/ICPR.2018.8545506. 632 Xiaofeng Liu, Bo Hu, Xiongchang Liu, Jun Lu, Jane You, and Lingsheng Kong. Energy-constrained 633 self-training for unsupervised domain adaptation, 2021. URL https://arxiv.org/abs/ 634 2101.00316. 635 Xiaofeng Liu, Chaehwa Yoo, Fangxu Xing, Hyejin Oh, Georges El Fakhri, Je-Won Kang, and 636 Jonghye Woo. Deep unsupervised domain adaptation: A review of recent advances and per-637 spectives, 2022. URL https://arxiv.org/abs/2208.07422. 638 639 Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I. Jordan. Conditional adversarial 640 domain adaptation, 2018. URL https://arxiv.org/abs/1705.10667. 641 Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In International Confer-642 ence on Learning Representations, 2019. URL https://openreview.net/forum?id= 643 Bkg6RiCqY7. 644 Guangting Mai, James W. Minett, and William S.-Y. Wang. Delta, theta, beta, and gamma 645 brain oscillations index levels of auditory sentence processing. NeuroImage, 133:516–528, 646 2016. ISSN 1053-8119. doi: https://doi.org/10.1016/j.neuroimage.2016.02.064. URL https: 647 //www.sciencedirect.com/science/article/pii/S1053811916001737.

- Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh. Domain adaptation: Learning bounds and algorithms. *CoRR*, abs/0902.3430, 2009. URL http://arxiv.org/abs/0902.3430.
- Yishay Mansour, Mehryar Mohri, and Afshin Rostamizadeh. Multiple source adaptation and the renyi divergence. CoRR, abs/1205.2628, 2012. URL http://arxiv.org/abs/1205. 2628.
- Stephanie Martin, Peter Brunner, Chris Holdgraf, Hans-Jochen Heinze, Nathan E. Crone, Jochem Reiger, Gerwin Schalk, Robert T. Knight, and Brian N. Pasley. Decoding spectrotemporal features of overt and covert speech from the human cortex. *Frontiers in Neuroengineering*, 7, May 2014. doi: 10.3389/fneng.2014.00014.
- Paul Mermelstein. Distance measures for speech recognition, psychological and instrumental.
 Pattern Recognition and Artificial Intelligence, pp. 374–388, 1976. URL https://api.semanticscholar.org/CorpusID:208101681.
- Daniel Moyer, Greg Ver Steeg, Chantal M. W. Tax, and Paul M. Thompson. Scanner invariant representations for diffusion mri harmonization. *Magnetic Resonance in Medicine*, 84(4):2174–2189, 2020. doi: https://doi.org/10.1002/mrm.28243. URL https://onlinelibrary.wiley.com/doi/abs/10.1002/mrm.28243.
- Fabian Pedregosa, Gael Varoquaux, Alexandre Gramfort, Vincent Michel, Bertrand Thirion, Olivier
 Grisel, Mathieu Blondel, Peter Prettenhofer, Ron Weiss, Vincent Dubourg, Jake Vanderplas,
 Alexandre Passos, David Cournapeau, Matthieu Brucher, Matthieu Perrot, Edouard Duchesnay,
 and Gilles Louppe. Scikit-learn: Machine learning in python. *Journal of Machine Learning Research*, 12, 01 2012.
- Vitória Piai, Ardi Roelofs, and Eric Maris. Oscillatory brain responses in spoken word production reflect lexical frequency and sentential constraint. *Neuropsychologia*, 53:146–156, 2014. ISSN 0028-3932. doi: https://doi.org/10.1016/j.neuropsychologia.2013.11.014. URL https://www.sciencedirect.com/science/article/pii/S0028393213004119.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya
 Sutskever. Learning transferable visual models from natural language supervision, 2021. URL https://arxiv.org/abs/2103.00020.
- Harsh Rangwani, Sumukh K Aithal, Mayank Mishra, Arihant Jain, and R. Venkatesh Babu. A closer
 look at smoothness in domain adversarial training, 2022. URL https://arxiv.org/abs/
 2206.08213.
- Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In Nassir Navab, Joachim Hornegger, William M. Wells, and Alejandro F. Frangi (eds.), *Medical Image Computing and Computer-Assisted Intervention – MICCAI* 2015, pp. 234–241, Cham, 2015. Springer International Publishing. ISBN 978-3-319-24574-4.
- Swami Sankaranarayanan, Yogesh Balaji, Carlos D. Castillo, and Rama Chellappa. Generate to adapt: Aligning domains using generative adversarial networks, 2018. URL https://arxiv.org/abs/1704.01705.
- Jan-Mathijs Schoffelen, Robert Oostenveld, Nietzsche Lam, Julia Udden, Annika Hultén, and Peter Hagoort. A 204-subject multimodal neuroimaging dataset to study language processing. *Scientific Data*, 6, 04 2019. doi: 10.1038/s41597-019-0020-y.
- Meredith Shafto, Lorraine Tyler, Marie Dixon, Jason Taylor, James Rowe, Rhodri Cusack, Andrew Calder, William Marslen-Wilson, John Duncan, Tim Dalgleish, Richard Henson, Carol Brayne, and Fiona Matthews. The cambridge centre for ageing and neuroscience (cam-can) study protocol: A cross-sectional, lifespan, multidisciplinary examination of healthy cognitive ageing. *BMC neurology*, 14:204, 10 2014. doi: 10.1186/s12883-014-0204-1.

694

 Hidemasa Takao, Naoto Hayashi, and Kuni Ohtomo. Effect of scanner in longitudinal studies of brain volume changes. *Journal of magnetic resonance imaging : JMRI*, 34:438–44, 08 2011. doi: 10.1002/jmri.22636.

702 703 704 705 706	Hidemasa Takao, Naoto Hayashi, and Kuni Ohtomo. Effects of study design in multi-scanner voxel- based morphometry studies. <i>NeuroImage</i> , 84:133–140, 2014. ISSN 1053-8119. doi: https: //doi.org/10.1016/j.neuroimage.2013.08.046. URL https://www.sciencedirect.com/ science/article/pii/S1053811913009099.
708 707 708 709 710 711 712	Jason R. Taylor, Nitin Williams, Rhodri Cusack, Tibor Auer, Meredith A. Shafto, Marie Dixon, Lorraine K. Tyler, Cam-CAN, and Richard N. Henson. The cambridge centre for ageing and neu- roscience (cam-can) data repository: Structural and functional mri, meg, and cognitive data from a cross-sectional adult lifespan sample. <i>NeuroImage</i> , 144:262–269, 2017. ISSN 1053-8119. doi: https://doi.org/10.1016/j.neuroimage.2015.09.018. URL https://www.sciencedirect. com/science/article/pii/S1053811915008150. Data Sharing Part II.
713 714	Eric Tzeng, Judy Hoffman, Trevor Darrell, and Kate Saenko. Simultaneous deep transfer across domains and tasks, 2015. URL https://arxiv.org/abs/1510.02192.
715 716 717	Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation, 2017. URL https://arxiv.org/abs/1702.05464.
718 719	Laurens van der Maaten and Geoffrey Hinton. Viualizing data using t-sne. Journal of Machine Learning Research, 9:2579–2605, 11 2008.
720 721 722	Agustin Vicente and Peter Langland-Hassan. <i>Inner Speech: New Voices</i> . Oxford University Press, 10 2018. ISBN 9780198796640.
723 724 725	Sarah Wandelt, David Bjanes, Kelsie Pejsa, Brian Lee, Charles Liu, and Richard Andersen. Representation of internal speech by single neurons in human supramarginal gyrus. <i>Nature Human Behaviour</i> , 8:1–14, 05 2024. doi: 10.1038/s41562-024-01867-y.
726 727 728 729	Fenqiang Zhao, Zhengwang Wu, Li Wang, Weili Lin, SHUNREN XIA, and Dinggang Shen. Har- monization of infant cortical thickness using surface-to-surface cycle-consistent adversarial net- works. In <i>Medical Image Computing and Computer Assisted Intervention - Conference Proceed- ings</i> , pp. 475–483. Springer Nature, oct 2019.
730 731 732	Yang Zou, Zhiding Yu, Xiaofeng Liu, B. V. K. Vijaya Kumar, and Jinsong Wang. Confidence regularized self-training, 2020. URL https://arxiv.org/abs/1908.09822.
733 734 735	A APPENDIX
736 737	A.1 Compute
738	As all experiments were carried out on a shared resource used by multiple research groups, a single
739	GPU was used in order to reduce wait times. Additionally, limitations related to specific configura-
740	tions during the period of this project's completion meant full-run experiments on the MEGalodon
741	architecture were not feasible to carry out. As it stands, even relying on subsets approximately 15%
742	1,600 hours of GPU compute in the completion of this study.
744 745	A.2 ADDITIONAL PREPROCESSING DETAILS
746	All work using the MEGalodon framework as a base applied standard procedures using the native
747	functionality of the MEG dataloaders from the PNPL python library ⁶ . A low-pass filter was applied
740	at 125 Hz and high-pass filter at 0.5 Hz in order to remove artifacts from muscle movements and

functionality of the MEG dataloaders from the PNPL python library⁶. A low-pass filter was applied at 125 Hz and high-pass filter at 0.5 Hz in order to remove artifacts from muscle movements and slow-drift, respectively. Additionally, a notch filter is applied at multiples of 50 Hz to account for possible line noise from the electric grid where the original recordings were taken. The signal is also downsampled to 250 Hz, taking care to avoid aliasing at frequencies up to the threshold set by the low-pass filter. Bad sensor channels are then detected with a variance threshold and replaced by interpolation from the nearest sensors (Jayalath et al., 2024). The Brainmagick framework has unique requirements in the way data is sliced, batched, and tracked and therefore is not able to rely

⁷⁵⁵

⁶https://github.com/neural-processing-lab/pnpl



Figure 5: The normalized distributions of the sex of the subjects from the MOUS (Schoffelen et al., 2019), Cam-CAN (Shafto et al., 2014; Taylor et al., 2017), and Gwilliams et al. (2022) datasets. The density plotted along the y-axis represents the proportion (i.e. relative frequency) of each category within its respective dataset.



Figure 6: The normalized distributions of participant ages from subsets taken over the MOUS (Schoffelen et al., 2019) and Cam-CAN (Shafto et al., 2014; Taylor et al., 2017) datasets for ex-periments done with the MEGalodon (Jayalath et al., 2024) base architecture. The density plotted along the y-axis represents the proportion (i.e. relative frequency) of each category within its re-spective dataset. The mean age calculated over the two datasets is displayed by the dotted line. The age-balanced subsets were created by randomly selecting subjects from the overlap of the two whole dataset age distributions. The random subsets include subjects taken at random from the entire dis-tributions. In both cases the count of Male and Female participants is balanced with a tolerance of 1 subject in either direction.



Figure 7: The normalized distributions of participant sex from subsets taken over the MOUS (Schof-felen et al., 2019) and Cam-CAN (Shafto et al., 2014; Taylor et al., 2017) datasets for experiments done with the MEGalodon (Jayalath et al., 2024) base architecture. The density plotted along the yaxis represents the proportion (i.e. relative frequency) of each category within its respective dataset. The age-balanced subsets were created by randomly selecting subjects from the overlap of the two whole dataset age distributions. The random subsets include subjects taken at random from the entire distributions. In both cases the count of Male and Female participants is balanced with a tolerance of 1 subject in either direction.



Figure 8: The normalized distributions of participant ages and participant sexes from subsets taken
over the MOUS (Schoffelen et al., 2019) and Gwilliams et al. (2022) datasets for experiments done
with the Brainmagick (Défossez et al., 2023) base architecture. The density plotted along the y-axis
represents the proportion (i.e. relative frequency) of each category within its respective dataset. The
mean age calculated over the two datasets is displayed by the dotted line. Subsets were constructed
with the intention of mimicking the characteristics of the full distributions from which they were
sampled.

957 958

959

A.4 ADDITIONAL BRAINMAGICK DETAILS

960Défossez et al. (2023)'s architecture is composed generally of a speech module, a brain module, and961a contrastive loss. Wav2Vec 2.0 (Baevski et al., 2020), a self-supervised model trained on audio962alone, is used for the speech module as they find it best represents the latent representations of963speech sounds (Défossez et al., 2023). The brain module is constructed sequentially from a spatial964attention layer over the MEG (or EEG) sensors, a participant (1 × 1 convolution) layer, and a set of965convolutional blocks (Défossez et al., 2023).

965 966 967

A.4.1 SPATIAL ATTENTION

The spatial attention layer helps select the most salient sensors from the layouts used by different studies during collection when remapping the MEG data into a shared channel dimension. The design works by first projecting the three-dimensional sensor locations (i.e. input channels), *i*, to a two-dimensional plane. This is done using a function from the MNE (Gramfort et al., 2013) Python library that leverages a device-dependent surface meant to preserve channel distances. These two-

...

dimensional positions (x_i, y_i) are then normalized to [0, 1] and for each output channel, j, a function a_j over $[0, 1]^2$ is learnt. This function is parameterized in the Fourier space as $z_j \in \mathbb{C}^{K \times K}$ with K= 32 harmonics along each axis, giving the full function definition as

975 976 977

978

979

$$a_j(x,y) = \sum_{k=1}^{K} \sum_{\ell=1}^{K} \operatorname{Re}(z_j^{(k,\ell)}) \cos(2\pi(kx+\ell y)) + \operatorname{Im}(z_j^{(k,\ell)}) \sin(2\pi(kx+\ell y)).$$
(1)

The final weights over the input sensors are found by taking the softmax of the function a_j evaluated at the sensor locations (x_i, y_i) :

984

985

986

987

$$\forall j \in \{1, \dots, D_1\}, \, \mathsf{SA}(X)^{(j)} = \frac{1}{\sum_{i=1}^C e^{a_j(x_i, y_i)}} \left(\sum_{i=1}^C e^{a_j(x_i, y_i)} X^{(i)}\right) \tag{2}$$

where SA is the spatial attention (Défossez et al., 2023). Because a_j is periodic in practice, (x, y) are scaled down and a spatial dropout is applied by sampling a location and removing each sensor within a specified distance from the softmax.

988 989 A.4.2 CLIP Loss

990 Défossez et al. (2023)'s Brainmagick architecture uses a multi-modal CLIP (originally Contrastive 991 Language-Image Pre-Training) loss (Radford et al., 2021). Commonly, a regression loss is used in 992 the supervised training of decoders to predict latent representations of speech known to be relevant 993 to the brain - in many cases the Mel spectrogram due to its similarity to how sound is represented in 994 the cochlea (Mermelstein, 1976). The problem with regression objectives in this context, however, is that they rely on the assumption that the dimensions of the Mel spectrogram are all scaled correctly 995 and equally important. In reality, some (e.g. very low) frequencies are irrelevant to speech and 996 can be differentiated by irregular orders of magnitude. The CLIP loss importantly does not aim 997 to maximally distinguish speech segments from one another but acts to relax the constraints of a 998 regression loss which may be tied too heavily to the above assumptions of relevancy, accuracy, 999 and scaling with respect to the representations from the speech module (Défossez et al., 2023). 1000 Given a brain recording segment X and the representation of the corresponding speech sound $Y \in$ 1001 $\mathbb{R}^{F \times T}$, N-1 negative samples $Y_i \in \{1, \dots, N-1\}$ are taken over the dataset and a positive 1002 sample is added as $\tilde{Y}_N = Y$. The training objective therefore becomes predicting the probability 1003 $\forall j \in \{1, \dots, N\}, p_j = \mathbb{P} \left| \tilde{Y}_j = Y \right|$ such that the model \mathbf{f}_{clip} maps X to a latent representation 1004 $Z = \mathbf{f}_{clin}(X) \in \mathbb{R}^{F \times T}$. We can approximate the objective by taking the softmax of the dot product 1005 of Z and the candidate speech representations Y_j :

$$\hat{p}_j = \frac{e^{\langle Z, \tilde{Y}_j \rangle}}{\sum_{i'=1}^N e^{\langle Z, \tilde{Y}_{j'} \rangle}},\tag{3}$$

1008 1009 1010

1007

1011 where $\langle \cdot, \cdot \rangle$ is the inner product over both dimensions of Z and \tilde{Y} Défossez et al. (2023). The CLIP 1012 loss is thus the cross-entropy between p_j and \hat{p}_j , simplifying to:

$$L_{\text{CLIP}}(p,\hat{p}) = -\log(\hat{p}_N) = -\langle Z, Y \rangle + \log\left(\sum_{j'=1}^N e^{\langle Z, \tilde{Y}_{j'} \rangle}\right)$$
(4)

under the assumption of a dataset large enough that the probability of sampling the same segment twice can be neglected (Défossez et al., 2023).

1020 A.4.3 CONFUSION LOSS

1022 The formal definition of the confusion loss as introduced by Tzeng et al. (2015) and used by Dinsdale et al. (2021) is given as:

1024

1019

1021

$$L_{\text{conf}}(X_u, d_u, \Theta_d; \Theta_{\text{repr}}) = -\frac{1}{S_u} \sum_{s=1}^{S_u} \sum_{k=1}^N \frac{1}{N} \log(p_{s,k})$$
(5)

1029 1030 Encoder Block Task Head 1031 Linear 1032 Conv 1033 Einsum 1034 Spatial Attention 1035 Softmax Conv MEG 1036 Conv Conv Inputs CLIP Probabilities Conv Domain Classifie Subject Layer Linear 1039 for $s \in [S]$ GLU Rel U 1040 Subject Block Dropout 1041 Block repeated for k = 11043

where only the parameters in Θ_{repr} are updated depending on the fixed value of Θ_d as indicated by $L_{\text{conf}}(X_u, d_u, \Theta_d; \Theta_{\text{repr}})$.

Figure 9: The Brainmagick architecture (Défossez et al., 2023) as modified for feature-level harmonization. No activation functions are used in the subject block. In the five repeating convolutional blocks, the first two convolutions use a residual skip connection, increasing dilation, BatchNorm layer, and GELU activation. The final convolution in these blocks is not residual and halves the number of channels with a GLU activation. The encoder block then applies two 1×1 convolutions with a GELU activation after the first. We use same domain classifier as Dinsdale et al. (2021), and use the CLIP network for the primary decoding task. The einsum refers to the tensor operation to calculate the normalized similarity scores between candidate and estimate segments.

1052 1053

1055

1044

1054 A.4.4 DATALOADERS

The original code from Défossez et al. (2023) creates custom classes which track the raw .way 1056 files of the simulus and MEG recordings in blocks chunked by seconds moving forward in time. 1057 For the purposes of the current study, we adopt the convention from the original paper and use a 1058 6 second minimum block size. Additionally, because a linguistic representation is used directly 1059 in the decoding step, the custom dataclasses also enforce that across all individual subjects the train, validation, and test segments are mutually exclusive with respect to the sentences presented as 1061 stimulus. This is maintained for the case where the splits are built from more than one dataset. We 1062 further modify this behavior to ensure that data within every batch is tracked with an identifier for its dataset of origin. This information is then extracted at train time to form the ground truth vectors 1063 for the domain classifier. Support for the selection of specific subjects, enabling study of subsets of 1064 any size and construction, was also added during the completion of the present study.

1066 1067 1068

A.5 ADDITIONAL MEGALODON DETAILS

Jayalath et al. (2024) define three pretext tasks for speech decoding: band prediction, phase shift 1069 prediction, and amplitude scale prediction. The band prediction task randomly selects and applies a 1070 band-stop filter to one of the frequency bands typically associated with brain activity: Delta (0.1-4)1071 Hz), Theta (4-8 Hz), Alpha (8-12 Hz), Beta (12-30 Hz), Gamma (30-70 Hz), and High Gamma (¿70 1072 Hz) (Giraud & Poeppel, 2012; Piai et al., 2014; Mai et al., 2016). The goal is to then predict the frequency band which was rejected. The phase shift prediction task is similar in nature: a discrete 1074 uniform random phase shift is applied to a uniform randomly selected proportion of the MEG sen-1075 sors, with the goal of predicting which phase shift was applied (a discrete number of possible values are used in order to reduce the difficulty of the task by treating it as a multi-class problem) (Jayalath et al., 2024). The use of random sensors and uniform random selection is meant to mitigate the 1077 1078 effect of variance in sensor placement between studies by ensuring the differences between any two regions of the brain are represented. Finally, the amplitude scale prediction task selects a random 1079 proportion of the sensors and applies a discrete random amplitude scaling coefficient to the signal with the objective of predicting the scaling factor (Jayalath et al., 2024). The intention behind this task is to learn representations encoding relative sensor amplitude differences. These pretext tasks are used to pre-train the backbone, a dataset-conditional layer and encoder block, of the architecture (see fig 10) before being swapped out for the speech decoding tasks in the fine-tuning stage. Additionally, subject conditioning (via subject embeddings, in contrast with Défossez et al. (2023)) is applied at the bottleneck of the encoder block before the pre-text or fine-tuning task heads.



Figure 10: The MEGalodon architecture as modified for feature-level harmonization. Layers treated as part of the encoder block, task head, or domain classifier are respectively shown in yellow, purple, and green. All weights updated during pre-training are shown in blue. Weights trainable during fine-tuning are in red, with the addition of those in the encoder block in the case of deep fine-tuning.

1109

1111

1086

1110 A.5.1 DATALOADERS

The MEGalodon architecture is already capable of supporting multiple datasets during a single train-1112 ing run, and does this by leveraging the MultiDataloader class from the PNPL python library⁷. Under 1113 the original implementation, one batch from each dataset is returned in alternating fashion during 1114 training. Adversarial harmonization, however, is better served when every dataset is represented by 1115 at least one data-point in each batch. We accomplish this by creating a custom dataloader class, 1116 ComboLoader, able to take a list of other dataloaders and return batches in the form of a tuple 1117 containing a single batch from each of the original loaders. At train time, random slices, the sum 1118 of which is equivalent to the original batch size, are then taken from each batch in this tuple and 1119 processed. Upon aggregating, the effect becomes equivalent to that of a batch with the originally specified size but having a random mix of the data from every domain. The ground truth targets 1120 for the domain classifier are also calculated at train time using the lengths of the randomly gener-1121 ated batch slices. This approach allows for the implementation of adversarial harmonization with 1122 minimal changes to the underlying network architecture. The PNPL datasets additionally support 1123 returning metadata alongside each MEG recording, and this feature was used to extract the associ-1124 ated participant IDs which could then be used to retrieve the correct ages when creating the target 1125 vectors for age-based harmonization. 1126

1120

1133

1128 A.6 SOFTWARE CHALLENGES

While the use of popular deep learning libraries such as Pytorch and Lightning has many advantages, it can also lead to unexpected roadblocks when attempting to implement behaviors outside the expected scope of their standard workflows. This was the case for the present work when attempting to build out the functionality for adversarial harmonization. As a reminder, the adversarial phase of

⁷https://github.com/neural-processing-lab/pnpl

1134 the harmonization framework we implement is composed of three steps: (1) optimizing the encoder 1135 block and task head for the task of interest, (2) optimizing the domain classifier to maximize its 1136 ability to identify the target bias, and (3) optimizing the encoder block to erase any signal related 1137 to the target bias from its features. This means that every training step during this phase contains 1138 three backwards passes and three steps by different optimizers. However, the same initial set of features produced by the encoder block is used across all of these functions. Inevitably, this leads 1139 to clashes in the computational graph as the parameters of the encoder block are updated after the 1140 first optimizer's step call, but the same feature vector used in the third step is still associated with 1141 the previous version of those parameters. In older versions of Pytorch, setting the retain_graph 1142 parameter to True when calling the backward pass successfully navigated this issue. However, this 1143 behavior was not maintained for manual optimization in Lightning. Instead, we were compelled to 1144 implement a work-around which involved re-writing the forward pass for all layers of the encoder 1145 block such that the parameters of that layer are cloned and passed to its respective Pytorch Functional 1146 variant alongside the input. It should be noted that slight differences in the underlying construction 1147 of these layers from their named counterparts does introduce a numerical variance from that of the 1148 original models, however, it is on a order of magnitude small enough that it did not relevantly impact performance. The versions of all tensors related to one training step were examined to ensure that 1149 optimization continued to perform as expected when applying this procedure. 1150

1151 1152

1153

1173 1174 A.7 ADDITIONAL RESULTS

The subset results for the augmented Brainmagick architecture are shown in Table 4. We conduct a 1154 one-sided independent samples t-test using the subset results collected across three seeds. We find 1155 that the effect of adversarial harmonization on top-10 accuracy is statistically significant (p < 0.05) 1156 when evaluating on both the Gwilliams Gwilliams et al. (2022) test split (p = 0.0021) and the 1157 MOUS Schoffelen et al. (2019) test split (p = 0.0358). This is again demonstrated when training 1158 over the full datasets as seen in Table 2 and we include the Top-1 accuracy results here in Table 5 for 1159 completeness. As in the Top-10 case, the effect is statistically significant (p < 0.05) when evaluating 1160 on both the Gwilliams Gwilliams et al. (2022) test split (p = 0.0163) and the MOUS Schoffelen et al. 1161 (2019) test split (p = 0.0141). These results clearly show the ability of adversarial harmonization 1162 to enhance deep-learning architectures for cross-dataset generalization of MEG speech decoding 1163 where they might otherwise be unable to do so.

Subset Results		Top-10 Accuracy	
Method	Training Data	Gwilliams	MOUS
Control	Gwilliams + MOUS	$65.9\%\pm0.2$	$57.7\% \pm 0.5$
Harmonized	Gwilliams + MOUS	67.8% ± 0.2	$59.6\% \pm 0.0$

Table 4: As in the full-run case, we report Top-10 segment-level accuracy. The best performance recorded over the test split of each dataset subset is marked in bold. Confidence intervals are calculated over 3 seeds.

175	Full-Run Results		Top-1 Accuracy	
176	Method	Training Data	Gwilliams	MOUS
78	Control (Official repo)	Gwilliams	41.2%, 41.3%*	-
70	Control (Official repo)	MOUS	-	40.4%, 36.8%*
00	Control (Our implementation)	Gwilliams	69.8%	-
50	Control (Our implementation)	MOUS	-	37.8%
81	Pre-trained on MOUS	Gwilliams	39.4%	-
2	Pre-trained on Gwilliams	MOUS	-	36.9%
3	Control	Gwilliams + MOUS	$39.2\%\pm0.5$	$36.6\%\pm0.5$
34	Harmonized	Gwilliams + MOUS	$\textbf{41.4\%}\pm0.3$	$38.8\% \pm 0.2$

1185

Table 5: Top-1 segment-level accuracy with confidence intervals calculated over 3 seeds. Results as reported in the original study are denoted by a single asterisk (*). The best performance recorded over each validation dataset is marked in bold.







Below in figures 14 and 15 we show that, holding all else equal, beginning the harmonization phase
later into training does not mitigate the tendency of the task loss to diverge. The following plots do,
however, demonstrate that after an initial peak following the start of harmonization, the task loss
once again begins to trend downwards. Note that convergence (as determined by early stopping)

is not reached by the control even within 400 epochs, which is why we opt for an epoch-based scheduling of the warm-up phase in the case of MEGalodon.

Figures 16 and 17 demonstrate the case of training on smaller batch sizes. Overall loss is reduced, but the model still fails to reach convergence within 200 epochs and thus beginning harmonization during this time still leads to divergence.





25





A.10 HYPERPARAMETERS

1406	Parameter	Value	Final Validation Loss
1407	Harmonization phase start	Epoch 25	7.61
1409	Harmonization phase start	Epoch 100	5.76
1410	Alpha	1.0	5.76
1411	Alpha	0.33	7.30
1412	Alpha Optimizer	0.25 Adam	5.65 7.00
1413	Harmonization phase LR	0.000001	5.76
1414	Harmonization phase LR	0.000005	5.88
1415	Harmonization phase LR	0.000066	5.98

Table 7: Hyperparameter testing of the augmented MEGalodon architecture carried out over subsets of the MOUS and Cam-CAN datasets. Tested parameters are listed, with all other values for that run held constant. Validation loss is reported at epoch 200. The final configuration of the model used for running experiments is an alpha of 0.25, beta of 1 (choice of beta value was negligible), harmo-nization phase start of 100, and harmonization learning rate of 0.00002 for the task and classifier optimizers and 0.00001 for the adversarial optimizer. The choice to set the adversarial learning rate to half that of the others came recommended by Dinsdale et al. (2021) to increase training stability.

A.11 LIMITATIONS

The present study was limited by time and resource constraints which ultimately meant results could not be collected across multiple seeds in all cases. Additionally, testing over the full datasets was not carried out for the experiments using the MEGalodon Jayalath et al. (2024) base architecture. Given the flexibility of the chosen harmonization framework, the present study would have benefited from exploring its capacity to combine more than two datasets at a time. An initial look at pooling three datasets for pre-training was done but not investigated further within the scope of this work. A significant bug in the code was discovered relatively late into the project which forced the experimental results collected up to that point to have to be discarded and re-collected. This setback meant that more extensive testing of the kind discussed above was infeasible. While this is regrettable, bugs in large and complex codebases, particularly those that build on other's publicly available code, can be commonplace. It is important to us to have caught the bug and be able to present accurate results, rather that allow it to remain undiscovered. An extension of the framework to enable harmonization of datasets with skewed demographic biases, such as the age distributions of MOUS Schoffelen et al. (2019) and Cam-CAN Shafto et al. (2014); Taylor et al. (2017), is noted in Dinsdale et al. Dinsdale et al. (2021). In this variant, the datasets are trained on in full, but harmonization is only carried out using subjects from the overlapping area of the distributions. This extension was implemented but was not able to be properly tested at the current time.