
MAEEG: Masked Auto-encoder for EEG Representation Learning

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Decoding information from bio-signals such as EEG, using machine learning
2 has been a challenge due to the small data-sets and difficulty to obtain labels.
3 We propose a reconstruction-based self-supervised learning model, the masked
4 auto-encoder for EEG (MAEEG), for learning EEG representations by learning
5 to reconstruct the masked EEG features using a transformer architecture. We
6 found that MAEEG can learn representations that significantly improve sleep stage
7 classification ($\sim 5\%$ accuracy increase) when only a small number of labels are
8 given. We also found that input sample lengths and different ways of masking
9 during reconstruction-based SSL pretraining have a huge effect on downstream
10 model performance. Specifically, learning to reconstruct a larger proportion and
11 more concentrated masked signal results in better performance on sleep classifica-
12 tion. Our findings provide insight into how reconstruction-based SSL could help
13 representation learning for EEG.

14 1 Introduction

15 Decoding information from Electroencephalography (EEG) and other bio-signal modalities has
16 enabled health-related clinical applications, such as sleep pattern detection [Ghassemi et al., 2018].
17 Recently, machine learning models have been effectively applied on classifying EEG signals, mostly
18 via supervised learning [Lawhern et al., 2018]. However, such methods rely on a large number
19 of labels for bio-signal data, which are usually difficult and expensive to obtain. One way to
20 improve accuracy when decoding EEG signals using only a small number of labels is to learn useful
21 representations from unlabeled data via self-supervised learning (SSL).

22 SSL has been widely explored in the field of computer vision and speech processing. The idea of SSL
23 is to boost the performance of a downstream task by learning meaningful representations through
24 an SSL task. One way to learn representations is through masking and reconstructing the input
25 signals (i.e., reconstruction-based SSL). For example, the representations of speech learned from the
26 Wav2Vec series of SSL models significantly improve the performance of various downstream tasks,
27 such as automatic speech recognition and speech emotion recognition [Baevski et al., 2020, Mitra
28 et al., 2022]. Given the success of reconstruction-based SSL methods in other domains, it is desirable
29 to learn whether and how such methods could be adapted to physiological time-series data, such as
30 EEG.

31 In this study, we explore representation learning using reconstruction-based SSL on EEG data. We
32 propose an SSL model, masked auto-encoder for learning EEG representations (MAEEG), which
33 can learn EEG representations by reconstructing the raw signal from masked features. We found
34 that MAEEG pretraining learns meaningful EEG representations, which yield better performance
35 on sleep stage classification. We further explore how masking may affect SSL and downstream task
36 performance, and found that in general, a higher probability and more concentrated masking yields
37 better task performance.

38 2 Background and Related Work

39 SSL methods for time-series data have been widely applied in domains such as natural language
40 and speech processing. Recently, studies have also explored various SSL methods for bio-signal
41 data, such as EEG. For example, combining various ways of EEG data augmentation and contrastive
42 learning, Mohsenvand et al. [2020] showed improved performance on several EEG classification tasks.
43 On the other hand, Banville et al. [2021] proposed three SSL methods aiming to learn representations
44 by discriminating the temporal relationships between the EEG samples. These studies showed that a
45 SSL-pretrained encoder could produce features that outperformed supervised learning when there is
46 only a small number of labeled data.

47 Another type of SSL methods for time-series data learn representations through masking and
48 reconstructing features, and we refer to them as reconstruction-based SSL in the current study.
49 Reconstruction-based SSL usually contains several stages: (1) The raw input signal is first encoded
50 as features using a convolutional encoder; (2) part of the features are “masked” by setting the features
51 to a certain value; (3) the masked features are sent to a second encoder, usually a transformer encoder,
52 which aims to reconstruct the masked part of the features using the unmasked features as context;
53 (4) finally, contrastive loss, or reconstruction loss is calculated to optimize the SSL task. Inspired
54 by the Wav2Vec2 model, Kostas et al. [2021] proposed BENDR, a reconstruction-based SSL with
55 contrastive learning for processing EEG signals. While the idea of using massive data (i.e., the TUH
56 dataset, [Obeid and Picone, 2016]) to learn EEG representations using BENDR seems to improve
57 some downstream tasks, the results are not consistent across different data-sets/tasks. This could
58 be due to the inconsistency of data-sets and model settings between pretraining and downstream
59 stages. Overall, it remains unclear if their reconstruction-based SSL methods could learn useful
60 representations for EEG signals.

61 In this study, we examine whether reconstruction-based SSL methods are effective for learning EEG
62 representations, and how different ways of masking would affect representation learning. Specifically,
63 we propose an SSL model, MAEEG, inspired by both BENDR and the MAE model proposed in
64 computer vision [He et al., 2022], to compare with BENDR and to broaden our understanding of
65 reconstruction-based SSL for EEG.

66 3 Method

67 3.1 Sleep EEG data-set

68 To examine how reconstruction-based SSL may help learn useful representations for classifying sleep
69 stages from EEG signals, we use the data-set provided from 2018 PhysioNet Challenge [Ghassemi
70 et al., 2018]. The data-set contains overnight polysomnography of 994 subjects (i.e., around 7,000
71 hours of data, 65% male, mean age: $55y \pm 14.4$), which was split into a training set (596 subjects),
72 a validation set (198 subjects) and a testing set (200 subjects) as in Ghassemi et al. [2018]. Each
73 polysomnogram contains simultaneous recordings of 6-channel EEG, ECG and respiratory signals
74 during sleep and two sets of labels: 5 sleep stages and respiratory-effort related arousals. The labels
75 were assigned by a trained sleep technician for each 30-sec non-overlapping window. Here, we
76 focused on using the 6-channel EEG signals to classify 5 stages of sleeping (i.e., wake, N1, N2, N3,
77 REM). All the models and experiments presented in this section are built and run on Pytorch 1.8.1.

78 3.2 Reconstruction-based SSL models

79 **BENDR.** The model architecture of BENDR is shown in Figure 1A. First, the 6-channel 100Hz
80 raw EEG signals are the input to the 6 convolutional layers, which encode the raw EEG into
81 64-dimensional convolved features (t_i , $\sim 1.05\text{Hz}$). Dropout and GroupNorm are added to each
82 convolutional layer and GELU is the output activation function. A mask I_m is generated to “mask out”
83 part of the features t_i where $i \in I_m$ (i.e., set the activation to Gaussian noise $M \sim \mathcal{N}(0, \sqrt{1/64})$).
84 The masked features q_i is then sent to an 8-layer transformer encoder and encoded as 192-dimensional
85 output features. The positions are encoded using a convolutional layer, which acts as relative positional
86 embedding. Finally, a final convolutional layer maps the features back to 64-dimensional contextual
87 features c_i . Contrastive loss is calculated by comparing c_i with the unmasked features t_i .

88 **MAEEG.** The MAEEG model has a similar architecture to BENDR (Figure 1B), but has two
 89 additional layers to map the transformer output back to the raw EEG dimensions. A reconstruction
 90 loss is calculated by comparing the reconstructed EEG (\hat{x}) and input EEG (x) signals as $1 - \frac{\hat{x} \cdot x}{\|\hat{x}\| \|x\|}$
 91 (Figure 1 B). The key difference between BENDR and MAEEG is that instead of using contrastive
 92 learning, MAEEG learns representations by minimizing the reconstruction loss.

93 3.3 Features and Analyses

94 After SSL pretraining, we examined downstream task performances by adding a linear classification
 95 layer on top of either of the two features: the convolved features t_i vs. the contextual features c_i ,
 96 under either of the two settings: freeze vs. fine-tune the encoder (See Appendix A.1, Figure A.1). We
 97 also investigated how masking in reconstruction-based SSL could affect EEG representation learning
 98 by modifying the masking probability and lengths (See Appendix A.2).

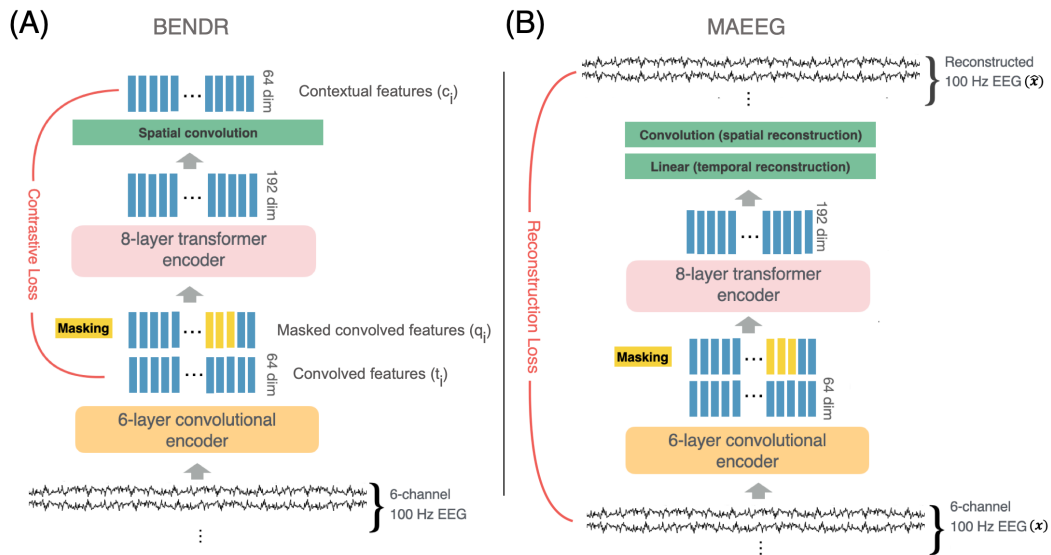


Figure 1: **Comparison between two reconstruction-based SSL model architecture.** **A.** BENDR learns representations through contrastive loss between convolved and contextual features, while **B.** MAEEG learns representations through reconstruction loss between input and output EEG signals.

99 4 Results

100 4.1 MAEEG learned useful representations for classifying sleep stages

101 We first examine convolved features learned solely by SSL pretraining by freezing the convolutional
 102 encoders when performing the sleep stage classification task. We found that representations pretrained
 103 by MAEEG (both 30s and 100s samples input) significantly outperformed the representations trained
 104 by BENDR and the supervised baseline model across different amount of subjects for training (Figure
 105 2A). This suggests that representations learned from MAEEG pretraining do capture features that are
 106 useful for classifying sleeping stages.

107 We next examined the model performance using contextual features when fine-tuning the encoders
 108 during the sleep stage classification task. We found that when only given one subject’s labels
 109 (containing all sleep stages), MAEEG pretrained on 100s (MAEEG-100s) outperformed other models,
 110 while BENDR pretrained on 30s (BENDR-30s) performed the best on 10 subjects’ and 50 subjects’
 111 labels (Figure 2B). The best classification performance (90% accuracy) was from MAEEG-100s.
 112 Interestingly, MAEEG-100s performed significantly better compared to MAEEG-30s, and an opposite
 113 effect on BENDR was observed - BENDR-30s performed significantly better compared to BENDR-
 114 100s, suggesting that representations learned from the two reconstruction-based SSL models are

115 distinct and the sample length for pretraining has a big effect on downstream performance. To
 116 understand how the pretraining differs under different experimental settings, we further examined the
 117 attention maps learned during pretraining. We found that models learned to attend distant features
 118 performed better than models learned to attend neighboring features (See Appendix A.3, FigureA.3).
 119 Also, the advantage of MAEEG-pretrained convolved features did not reflect in model performance
 120 when the encoders were fine-tuned, suggesting that while MAEEG may learn some useful features
 121 to discriminate sleep stages through SSL pretraining, such pretraining may also hinder the models’
 122 ability to learn at the same time.

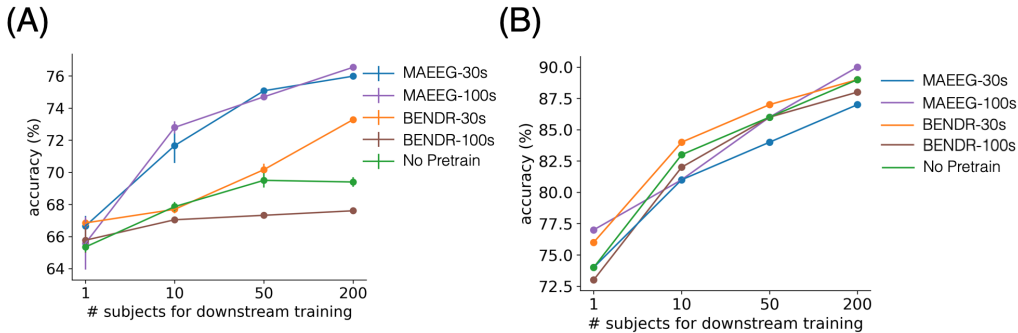


Figure 2: **Model comparison on classifying sleep stages A.** Classification performance with different numbers of downstream subjects using convolved features when the encoders were frozen. **B.** Classification performance using contextual features when the encoders were fine-tuned. Chance level was 20%. The error bars indicate standard deviation of 3 models with different weight initiation.

123 **4.2 Masking effects in MAEEG**

124 Next, we examined how different ways of masking during reconstruction-based SSL pretraining may
 125 affect representation learning and downstream classification performance using contextual feature c_i
 126 with finetuning, by varying the mask rate and the number of mask chunks (See Appendix A.2). We
 127 first examined the SSL loss during pretraining, and found that higher percentage of masking with less
 128 chunks (which yields longer masks) yielded higher SSL loss (i.e. more difficult to train). We then
 129 examined the downstream model performance and found that the model performances differ most
 130 when given the least number of labels - such as 0.1% (i.e. ~4 hours of labels) of downstream training
 131 data (Figure A.2B). With such small number of labels, we found that masking 75% of the tokens
 132 performed better overall compared to masking 50% and 25% of the tokens, and models with a single
 133 chunk performed better than models with multiple chunks. All models with MAEEG pretraining
 134 performed better than the models without pretraining (Figure A.2C). Finally, we examined whether
 135 such a trend (single chunk with longer mask performs better) preserved when we vary the mask
 136 span of a single mask chunk, and whether it generalizes to BENDR. We found a consistent pattern
 137 across both MAEEG and BENDR, that models pretrained with longer mask span performs better than
 138 models pretrained with shorter mask span when there only a small number of labels (Figure A.2D).
 139 This finding is likely due to the fact that a single and long mask during pretraining forced the model
 140 to learn a more difficult task and yield more useful representations for sleep stage classification.

141 **5 Conclusion**

142 Representation learning for EEG can be useful given the difficulty in collecting large amount of data
 143 and labels. In this study, we proposed a reconstruction-based SSL model, MAEEG, for learning
 144 EEG representations by reconstructing EEG signals using masked auto-encoder architecture. We
 145 tested our model on sleep stage classification and found that MAEEG can learn useful representations
 146 solely from the pre-text task which boost the downstream classification performance. We also found
 147 that different ways of masking and sample lengths selected during SSL pre-training can significantly
 148 affect the downstream classification performance. We encourage future studies to examine how
 149 reconstruction-based SSL may help representation learning on other time-series data for health.

150 References

- 151 Alexei Baevski, Yuhao Zhou, Abdelrahman Mohamed, and Michael Auli. wav2vec 2.0: A framework
152 for self-supervised learning of speech representations. *Advances in Neural Information Processing*
153 *Systems*, 33:12449–12460, 2020.
- 154 Hubert Banville, Omar Chehab, Aapo Hyvärinen, Denis-Alexander Engemann, and Alexandre
155 Gramfort. Uncovering the structure of clinical eeg signals with self-supervised learning. *Journal*
156 *of Neural Engineering*, 18(4):046020, 2021.
- 157 Mohammad M Ghassemi, Benjamin E Moody, Li-Wei H Lehman, Christopher Song, Qiao Li, Haoqi
158 Sun, Roger G Mark, M Brandon Westover, and Gari D Clifford. You snooze, you win: the
159 physionet/computing in cardiology challenge 2018. In *2018 Computing in Cardiology Conference*
160 *(CinC)*, volume 45, pages 1–4. IEEE, 2018.
- 161 Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked
162 autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on*
163 *Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.
- 164 Demetres Kostas, Stephane Aroca-Ouellette, and Frank Rudzicz. Bendr: using transformers and a
165 contrastive self-supervised learning task to learn from massive amounts of eeg data. *Frontiers in*
166 *Human Neuroscience*, page 253, 2021.
- 167 Vernon J Lawhern, Amelia J Solon, Nicholas R Waytowich, Stephen M Gordon, Chou P Hung, and
168 Brent J Lance. Eegnet: a compact convolutional neural network for eeg-based brain–computer
169 interfaces. *Journal of neural engineering*, 15(5):056013, 2018.
- 170 Vikramjit Mitra, Hsiang-Yun Sherry Chien, Vasudha Kowtha, Joseph Yitan Cheng, and Erdrin
171 Azemi. Speech emotion: Investigating model representations, multi-task learning and knowledge
172 distillation. *arXiv preprint arXiv:2207.03334*, 2022.
- 173 Mostafa Neo Mohsenvand, Mohammad Rasool Izadi, and Pattie Maes. Contrastive representation
174 learning for electroencephalogram classification. In *Machine Learning for Health*, pages 238–253.
175 PMLR, 2020.
- 176 Iyad Obeid and Joseph Picone. The temple university hospital eeg data corpus. *Frontiers in*
177 *neuroscience*, 10:196, 2016.

178 A Appendix

179 A.1 Features for Downstream classification

180 As mentioned, two features are learned by the SSL models - the convolved features, t_i , which are
181 the output of the convolutional encoder, and the contextual features, c_i , which are the output of
182 the transformer encoder. To build downstream classification model for classifying sleep stages, we
183 preserve most of the SSL architecture with the pretrained encoders, and simply add a fully-connected
184 classification layer which takes the SSL features as input and generates probabilities for predicting
185 the 5 sleep stages (Figure A.1). To better understand representations learned from different stages
186 of training (i.e., SSL pretraining with no labels and downstream supervised training with labels),
187 we conducted two types of downstream classification analyses: (1) To examine low-level features
188 t_i learned solely from the SSL pre-training, we use t_i while “freezing” the convolutional encoder
189 (meaning that the weights of the encoder are fixed) during downstream classification. Only the
190 classification layer gets trained in this condition (Figure A.1 A); (2) To examine the pre-trained
191 model as an initialization for supervised training, we use c_i while fine-tuning both convolutional and
192 transformer encoders and the classification layer during downstream classification (Figure A.1 B).
193 For (2), to determine which transformer layer contains the most useful SSL features for classification,
194 we first probed the features from each layer and found that layer 2 showed the best classification
195 performance out of all layers. We therefore extracted c_i from layer 2 as the SSL features to conduct
196 all following analyses. We conducted these analyses on MAEEG, BENDR and a baseline supervised
197 model which has the same architecture as other models but the encoders are randomly initialized,
198 i.e., no SSL pretraining prior to the downstream classification task. We applied analyses with various

199 number of subjects or percent of sessions to examine model performance given different amount of
200 labels.

201 **A.2 Masking effect analyses**

202 The masking during SSL pretraining determines the difficulty of the SSL task, and the representations
203 learned from such task. Therefore, it is critical to understand what could be an ideal way of masking
204 for learning EEG signals. In the original BENDR, they set a probability ($p=0.065$) for each token
205 to be the beginning of a mask with a length of 10 tokens. The problem with this is that for each
206 sequence, the mask can vary a lot - sometimes most of the tokens can be masked with overlap, and
207 sometimes there is no token being masked - which may cause the instability of representation learning.
208 Here, we tried to use a more systematic way to generate the mask, by determining the mask rate and
209 the number of mask chunks. These two factors would determine how long each mask looks like, as
210 illustrated in Figure A.2 A. We did two analyses:

- 211 1. Pretrained the models with 3 mask rates: 75%, 50% and 25%, and 3 numbers of mask
212 chunks: 1, 5, and 10.
- 213 2. Pretrained models with single mask chunk but vary the mask span to further confirm the
214 effects.

215 The resulting pretrained models were evaluated by the downstream performance when using the
216 learned features for classification. We conducted all the analyses on BENDR and MAEEG pretrained
217 on 100s samples because models pretrained on 30s samples did not have enough tokens for examining
218 the effect of these mask variations.

219 **A.3 Attention learned from SSL**

220 To understand how the EEG representation learned from reconstruction-based SSL help sleep classi-
221 fication, we visualized the attention from the transformer layer where features were extracted (i.e.
222 layer 2) after SSL pretraining but before extensive downstream supervised training (we slightly
223 tuned the models with 1 epoch with 1 subject's data just to get stable representations.) We extract
224 the token-by-token attention maps for the baseline supervised model and the models pretrained by
225 MAEEG and BENDR. We visualized the attention maps from a sample where the label is N2 (Figure
226 A.3). Compared to the attention from baseline models which is equally distributed as shown in
227 Figure A.3 A, MAEEG- and BENDR- pretrained models showed distinct attention patterns learned
228 from pretraining. Interestingly, we observed that MAEEG-100s and MAEEG-30s actually learned
229 very different attention patterns, with MAEEG-30s learning to attend locally to neighboring tokens
230 while MAEEG-100s learning to attend more distant tokens (Figure A.3 B). On the other hand,
231 BENDR-100s and BENDR-30s also learned very different attention patterns, with BENDR-100s
232 learning to attend locally and BENDR-30s learning a more global distributed attention. This suggests
233 that representation learning can be very different based not only on the model architecture, but
234 also on subtle SSL settings such as the length of pretraining samples. Moreover, MAEEG-30s and
235 BENDR-100s performed worse on the classification task compared to their counterparts (Figure 2B),
236 suggesting that local attention acquired from pretraining may not be desirable for the sleep stage
237 classification task, possibly due to the fact that this specific task relies on long-timescale information.
238

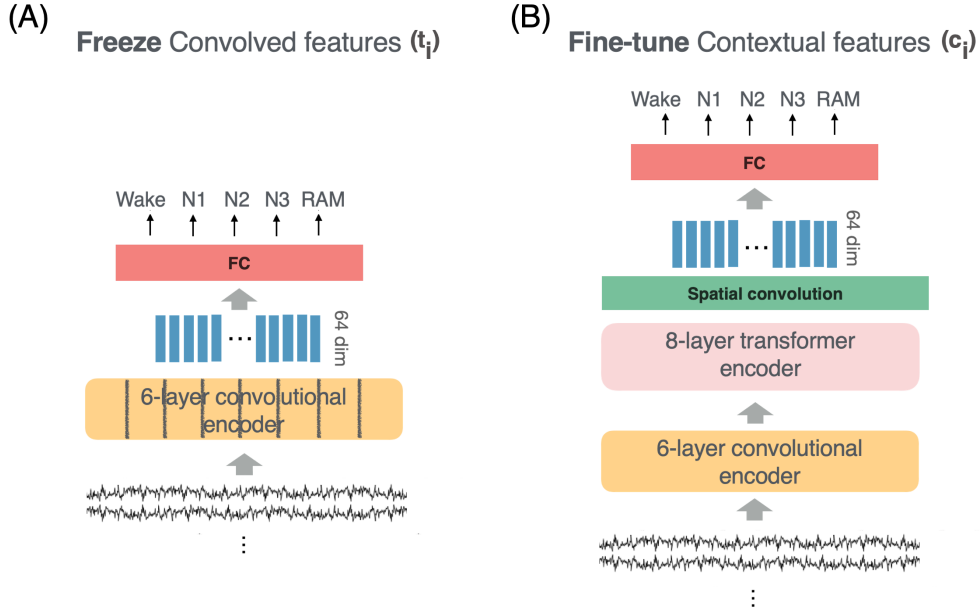


Figure A.1: **Downstream model architecture** **A.** Classification using convolved features t_i with encoder frozen to examine the representation learned solely by SSL without labels. **B.** Classification using contextual features c_i . The encoders are fine-tuned during downstream training to examine the pre-trained model as an initialization for supervised training.

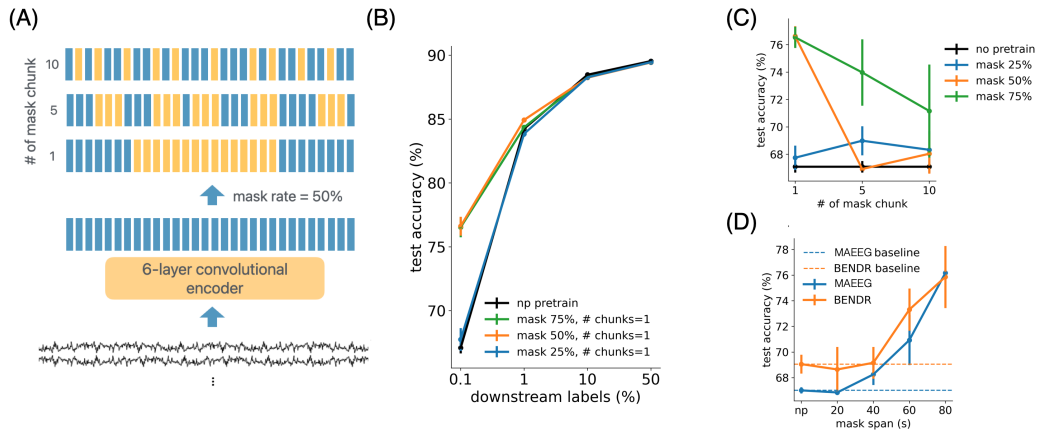


Figure A.2: **Masking effect on classification performance** **A.** Example of generating different masks for SSL pretraining by varying the mask rate and number of mask chunks. **B.** Classification performance using fine-tuned contextual features with various masking conditions and percentage of training sessions. **C.** Classification results using 0.1% of labels with different masking conditions. Overall, we found that masking 75% of the signals with 1 chunk results in the best performance. **D.** Classification results using 0.1% of labels when varying the masking span. We found that longer mask span yields better results for both MAEEG-100s and BENDR-100s. np = no pretrain

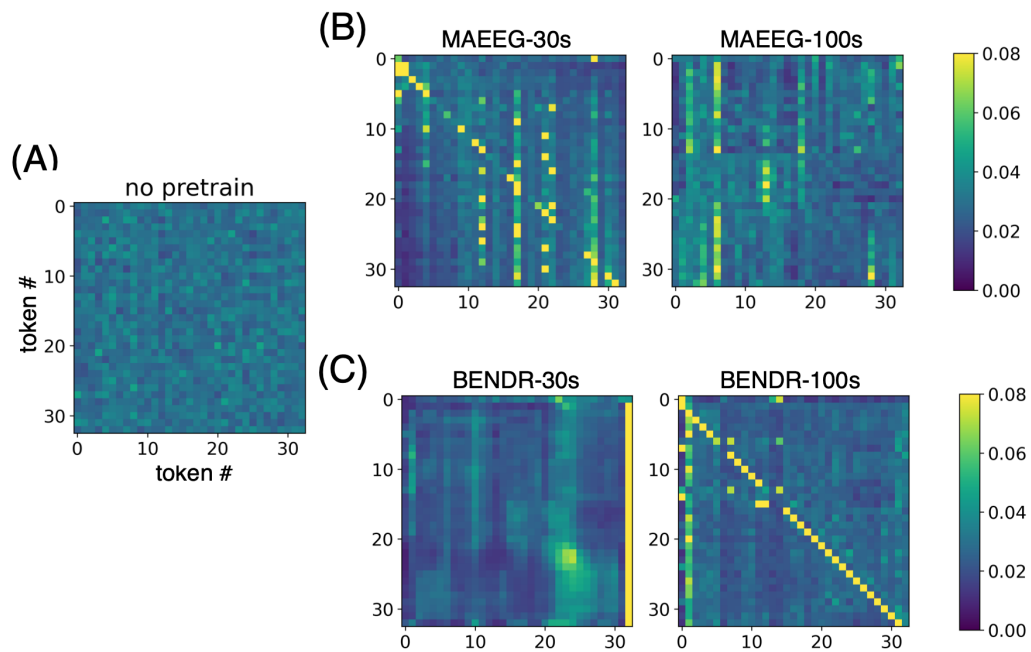


Figure A.3: **Visualizing attention after pretraining** **A.** Initialized attention. **B.** (left) attention learned from MAEEG-30s pretraining; (right) attention learned from MAEEG-100s pretraining **C.** (left) attention learned from BENDR-30s pretraining; (right) attention learned from BENDR-100s pretraining.